

BMJ Open

BMJ Open is committed to open peer review. As part of this commitment we make the peer review history of every article we publish publicly available.

When an article is published we post the peer reviewers' comments and the authors' responses online. We also post the versions of the paper that were used during peer review. These are the versions that the peer review comments apply to.

The versions of the paper that follow are the versions that were submitted during the peer review process. They are not the versions of record or the final published versions. They should not be cited or distributed as the published version of this manuscript.

BMJ Open is an open access journal and the full, final, typeset and author-corrected version of record of the manuscript is available on our site with no access controls, subscription charges or pay-per-view fees (<http://bmjopen.bmj.com>).

If you have any questions on BMJ Open's open peer review process please email info.bmjopen@bmj.com

BMJ Open

Biomarker discovery studies for patient stratification using machine learning analysis of omics data: a scoping review

Journal:	<i>BMJ Open</i>
Manuscript ID	bmjopen-2021-053674
Article Type:	Original research
Date Submitted by the Author:	20-May-2021
Complete List of Authors:	Glaab, Enrico; University of Luxembourg, Luxembourg Centre for Systems Biomedicine Rauschenberger, Armin; University of Luxembourg, Luxembourg Centre for Systems Biomedicine Banzi, Rita; Mario Negri Institute for Pharmacological Research, Center for Health Regulatory Policies Gerardi, Chiara; Mario Negri Institute for Pharmacological Research, Center for Health Regulatory Policies Garcia, Paula; ECRIN, European Clinical Research Infrastructure Network Demotes, Jacques; ECRIN, European Clinical Research Infrastructure Network
Keywords:	BIOTECHNOLOGY & BIOINFORMATICS, NATURAL SCIENCE DISCIPLINES, STATISTICS & RESEARCH METHODS

SCHOLARONE™
Manuscripts



I, the Submitting Author has the right to grant and does grant on behalf of all authors of the Work (as defined in the below author licence), an exclusive licence and/or a non-exclusive licence for contributions from authors who are: i) UK Crown employees; ii) where BMJ has agreed a CC-BY licence shall apply, and/or iii) in accordance with the terms applicable for US Federal Government officers or employees acting as part of their official duties; on a worldwide, perpetual, irrevocable, royalty-free basis to BMJ Publishing Group Ltd ("BMJ") its licensees and where the relevant Journal is co-owned by BMJ to the co-owners of the Journal, to publish the Work in this journal and any other BMJ products and to exploit all rights, as set out in our [licence](#).

The Submitting Author accepts and understands that any supply made under these terms is made by BMJ to the Submitting Author unless you are acting as an employee on behalf of your employer or a postgraduate student of an affiliated institution which is paying any applicable article publishing charge ("APC") for Open Access articles. Where the Submitting Author wishes to make the Work available on an Open Access basis (and intends to pay the relevant APC), the terms of reuse of such Open Access shall be governed by a Creative Commons licence – details of these licences and which [Creative Commons](#) licence will apply to this Work are set out in our licence referred to above.

Other than as permitted in any relevant BMJ Author's Self Archiving Policies, I confirm this Work has not been accepted for publication elsewhere, is not being considered for publication elsewhere and does not duplicate material already published. I confirm all authors consent to publication of this Work and authorise the granting of this licence.

Title

Biomarker discovery studies for patient stratification using machine learning analysis of omics data: a scoping review

Authors

Enrico Glaab^{1,*}, Armin Rauschenberger¹, Rita Banzi², Chiara Gerardi², Paula Garcia³, Jacques Demotes-Mainard³, and the PERMIT Group

Affiliations

¹Luxembourg Centre for Systems Biomedicine, University of Luxembourg, Campus Belval, Esch-sur-Alzette, Luxembourg.

²Center for Health Regulatory Policies, Istituto di Ricerche Farmacologiche Mario Negri, Milan, Italy.

³European Clinical Research Infrastructure Network (ECRIN), Paris, France

*Correspondence: enrico.glaab@uni.lu; 14, avenue du Rock'n'Roll, L-4361 Esch-sur-Alzette, Luxembourg; Phone: +352-466644 6186

Word count

Abstract: 291

Main text: 5320

Keywords

Biomarkers, Machine Learning, Review, Omics, Stratification

Abstract

Objective: To review biomarker discovery studies using omics data for patient stratification which led to clinically validated FDA-cleared tests or laboratory developed tests, in order to identify common characteristics and derive recommendations for future biomarker projects.

Design: Scoping review.

Methods: We searched PubMed, EMBASE and Web of Science to obtain a comprehensive list of biomedical literature articles describing clinically validated biomarker signatures for patient stratification, derived using statistical learning approaches. All documents were screened to retain only peer-reviewed research articles, review articles, or opinion articles, covering supervised and unsupervised machine learning applications for omics-based patient stratification. Two reviewers independently confirmed the eligibility. Disagreements were solved by consensus. We focused the final analysis on omics-based biomarkers which achieved the highest level of validation, i.e., clinical approval of the developed molecular signature as a laboratory developed test or FDA approved tests.

Results: Overall, 234 articles fulfilled the eligibility criteria. The analysis of validated biomarker signatures identified multiple common methodological and practical features that may explain the successful test development and provide guidance for future biomarker projects. These include study design choices to ensure sufficient statistical power for model building and external testing, suitable combinations of non-targeted and targeted measurement technologies, the integration of prior biological knowledge, and the adequacy of statistical and machine learning methods for discovery and validation.

Conclusions: While most clinically validated biomarker models derived from omics data have been developed for biomedical decision making in oncology, first applications for non-cancer diseases highlight the potential of multivariate omics biomarker design for other complex disorders. Distinctive characteristics of prior success stories, such as early filtering and robust discovery approaches, continuous improvements in assay design and experimental measurement technology, and rigorous multi-cohort validation approaches, enable the derivation of specific recommendations for future studies.

Article Summary

- This scoping review provides a first integrative overview of biomarker discovery studies using omics data which led to clinically validated diagnostic and prognostic tools.
- It identified shared characteristics of successful omics-based biomarker studies, which may help to guide study design, discovery and validation methods for future projects.
- Recommendations derived from the review mainly provide guidance on optimizing the design of prospective studies, but also include suggestions for retrospective studies.
- The integration of diverse, multi-omics data sources for biomarker modeling is still an exception, but has the potential to provide more robust and reliable biomedical predictions.
- Further knowledge exchange among computational, experimental and clinical experts in the field is still needed to derive comprehensive guidelines for omics-based biomarker studies.

Introduction

Personalised medicine is a rapidly developing area in health care research and practice, which aims at providing more effective and safer therapies tailored to the individual patient, by exploiting subject-specific molecular, clinical and environmental data sources (Box 1).

One of the main tools used in personalised medicine and the focus of this survey is the machine learning (ML) analysis of omics profiling data to derive molecular biomarker signatures for disease- or drug-based patient stratification. The major goals behind ML-based omics biomarker development in this domain are to develop more reliable and robust tests for drug response prediction, early diagnosis, differential diagnosis or prognosis of the future clinical disease course. Omics-derived biomarker signatures may help to guide treatment decisions, and to focus therapies on the right populations in order to prevent overtreatment, increase success rates, and reduce costs. As a research and information tool, they may also enable a better monitoring of disease progression and treatment success, and guide new drug development and discovery. In contrast to classical single-molecule biomarker approaches, omics signatures have the potential to provide more sensitive, specific and robust predictions of disease-associated outcomes.

However, while biomarker discovery projects using omics data have already led to the successful development of clinically validated diagnostic and prognostic tests (1–10), many biomarker studies are not further pursued after early development stages or fail the translation in later clinical validation stages. Dedicated statistical and ML methodologies for omics biomarker discovery and validation have been published, as well as recommendations for the study design, implementation and reporting (11,12), but it is not clear which distinctive features and approaches characterize the studies that succeed in translating omics research findings into clinically validated tests.

As part of an ongoing EU project on “Personalised Medicine Trials” (PERMIT, <https://permit-eu.org>), funded within the H2020 framework, we have therefore investigated the current methodological practices for personalised medicine, covering ML approaches for omics-based patient stratification as one of the major focus areas. While a broader series of questions was established for the overall scoping review (16), for this manuscript, we focused our analysis on biomarker discovery studies that have led to successful, clinically validated FDA-cleared tests or laboratory developed tests (LDTs), to determine their shared and distinctive characteristics as compared to previous biomarker studies without clinical translation. In particular, we aimed to address the following more specific research questions:

- Which omics-derived biomarker discovery studies have previously led to clinically validated tests for patient stratification (LDTs or FDA-cleared tests)?
- What are the key characteristics that are shared by successful omics biomarker studies and that distinguish them from previously published biomarker studies which have not yet led to clinically validated tests?
- Which types of model building and validation methods have been used to develop clinically validated biomarker signatures, and what are the lessons learned and recommended workflows?
- Which recommendations and guidelines have previously been proposed to address common challenges in biomarker development using omics data?

1
2
3 These questions lend themselves for a scoping review, because omics-derived biomarker
4 development is still an evolving field, and a preliminary assessment of the potential scope and size
5 of the available biomedical literature on these topics is required as a first step for further follow-up
6 research. Therefore, the objective of this survey was to address the above questions by retrieving
7 and examining the current literature that describe biomarker discovery and validation studies using
8 omics data and ML approaches.
9

10 11 12 **Methods**

13
14 We conducted a scoping review following the methodological framework suggested by the Joanna
15 Briggs Institute (13). This framework consists of six stages: 1) identifying the research questions, 2)
16 identifying relevant studies, 3) study selection, 4) charting the data, 5) collating, summarising and
17 reporting results, and 6) consultation.
18

19 The scoping review approach was considered to be the most suitable to respond to the broad scope
20 and the evolving nature of the field. Compared to systematic reviews that aim to answer specific
21 questions, scoping reviews are used to present a general overview of the evidence pertaining to a
22 topic and they are useful to examine areas that are emerging, to clarify key concepts and identify
23 gaps (14,15). Before conducting the review, a study protocol was published on the online platform
24 Zenodo (16). Due to the iterative nature of scoping reviews, deviations from the protocol are
25 expected and duly reported when occurred. We used the PRISMA-ScR (Preferred Reporting Items
26 for Systematic reviews and Meta-Analyses extension for Scoping Reviews) checklist to report our
27 results (17) (Online supplementary file 1).
28
29
30

31 32 *Study identification*

33 Relevant studies and documents were identified, balancing feasibility with breadth and
34 comprehensiveness of searches. We searched PubMed, EMBASE and Web of Science (search
35 date: March 13, 2020) for articles describing supervised or unsupervised ML analyses for biomarker
36 discovery or personalised medicine, including both discovery and validation methods (see Fig. 1,
37 illustrating the keyword-based search strategy). We included journal publications and meeting
38 abstracts from international conferences and workshops. No other grey literature was included.
39 Online supplementary file 2 reports the detailed search strategies applied. We restricted inclusion to
40 reports published from January 2000 to April 2020 (covering also “online first” articles with official
41 publication date in the near future) in English, French, Spanish, Italian and German language.
42
43
44

45 46 *Eligibility criteria*

47 We included peer-reviewed methodology articles, review articles, opinion articles on supervised and
48 unsupervised ML methods for omics stratification and associated validation methods (addressing
49 accuracy, robustness, and clinical relevance). Only approaches tested on real-world biomedical data
50 were reviewed, while studies relying purely on simulated data were excluded. We also excluded
51 papers on biomarker methodologies without a demonstrated biomedical application, and those with
52 insufficient sample size (i.e., less than 50 samples per group used for the main conditions studied,
53 unless a dedicated power calculation was presented) or statistical validation (i.e., lack of clear
54 descriptions of cross-validation or external testing methodology, performance metrics and test
55 statistics). These exclusion criteria were not specified in the generic review protocol, but they were
56 agreed among the authors before starting the screening process.
57
58
59
60

To cover both data from original research papers and prior systematic reviews, we extracted information from three main article types: (1) applied research papers, (2) methodology articles with demonstrated applications, and (3) review articles on methods, applications and validation approaches.

Apart from these inclusion and exclusion criteria, for the final presentation of results, the statistical investigations covered all selected articles, whereas the detailed discussion of study characteristics in this paper focused mainly on the studies that led to clinically validated biomarker signatures tested on multiple cohorts with large sample sizes (i.e., studies using a power calculation to demonstrate the adequacy of the chosen sample sizes, or covering hundreds or thousands of samples per studied subject group).

Study selection

We exported the references retrieved from the searches into the online tool Rayyan (18). Duplicates were removed automatically using the reference manager Endnote X9 (Clarivate Analytics, Philadelphia, United States) and manually by the reviewers. One reviewer screened the titles and abstracts and retrieved full-text copies of potentially eligible reports for further assessment. Two reviewers independently confirmed the eligibility. Disagreements were solved by consensus.

Charting the data

We designed a data extraction form using Excel (Online supplementary file 3). General study characteristics extracted were for instance: authors, title, citation, type of publication (e.g., journal article, meeting abstract), study population and sample size (if applicable), methodology/study design, and outcome measures (if applicable). Specific items included key findings that relate to the review question; type of article/study (e.g., methodology, applied research, review– methods, review– applications, review–validation); generic ML domain (e.g., supervised/unsupervised); name of specific ML approach used.

To capture key findings that relate to the review question, relevant sentences were extracted from each reviewed article, and if needed, complemented by a brief explanatory remark by the reviewer.

The reviewers piloted the data extraction form using five records from the retrieved article collection. Two reviewers (EG, AR) working independently extracted the data from the included articles. In the case of disagreements, consensus was obtained by discussion.

In the final full-text review stage, the pre-selected articles were grouped by topic, categorizing articles into applied vs. methodological studies, supervised vs. unsupervised analyses, and assigning algorithm type identifiers to each article (review articles and papers on validation methodologies were considered as separate categories without a specific algorithm type assignment).

It was not within the remit of this scoping review to assess the methodological quality of individual studies included in the analysis.

Consultation exercise

The members of the PERMIT consortium, associated partners, and the PERMIT project Scientific Advisory Board discussed the preliminary findings of the scoping review in a 2-hour online workshop.

Patient and public involvement

The European Patients' Forum is a member of PERMIT project. Although not directly involved in the conduction of the scoping review, they received the draft review protocol for collecting comments and feedback.

Results

Study selection and general characteristics of reports

We retrieved 1164 abstracts from the literature search. After the removal of duplicates, we screened the remaining 1079 abstracts for eligibility. 502 records were excluded, while 577 abstracts were retained for the full-text assessment. We finally included 234 articles that passed all filtering criteria in the data extraction and analysis step (see flow chart in Fig. 2 and online supplementary file 2).

The full-text article review revealed that many studies did not meet the pre-defined inclusion criteria: 251 articles (44%) were removed because of an insufficient sample size, and 67 further articles (12%) were removed because they provided insufficient details on the validation results or methodology (see Fig. 2). This shows that the challenges of recruiting an adequate number of participants per study group or conducting sufficient omics profiling experiments for robust model building and validation are not met in a large proportion of omics biomarker studies, and that many of these studies lack adequate documentation for the study design and validation.

For the 234 selected articles, the majority (81%) rely entirely on an internal validation involving data from only a single cohort, whereas studies that use an external validation on an independent cohort are still underrepresented (only 12% of articles describe both an internal cross-validation and an external cohort validation, and an additional 7% include an external validation, but do not report internal cross-validation results). Moreover, when comparing the numbers of published studies involving different types of internal and external validations over different periods of time during the past 20 years, the relative proportion of studies including an external validation has only slightly increased in recent years (see Fig. 3).

Since a detailed discussion of all filtered articles is not within the scope of the present review, in the following, we focus on reviewing representative omics-based biomarkers studies which have achieved the highest level of validation, i.e., clinical approval of the developed molecular signature as an LDT or FDA approved test (see the overview of studies in Table 1). By investigating the shared and distinctive features of these successful studies, we also cover how they address common shortcomings and missing features of other reviewed studies, and summarize the lessons learned.

Success stories in omics-based biomarker signature development

Cancer approved omics-derived diagnostic tests (8 studies)

The first and most well-known omics-derived molecular test to receive FDA clearance was *MammaPrint*, a prognostic signature using the RNA expression activity of 70 genes to estimate the risk for distant tumor metastasis and recurrence in early-stage breast cancer patients (1,19–23). This test had been developed at the Netherlands Cancer Institute, using DNA microarray analysis to investigate primary breast tumors of 117 patients. Supervised ML was applied to the resulting data to identify a gene signature that was highly predictive of a short interval to distant metastases in lymph node negative patients (19).

1
2
3 A distinctive feature of the development approach behind this signature in comparison to other
4 reviewed studies was the multi-stage filtering and cross-validation strategy used in the initial
5 discovery study, which may explain the repeated confirmation of the signature in later validation
6 studies (1,20–23). From 25k genes represented on the DNA microarrays, only those significantly
7 regulated in more than 3 tumors out of the subset of 78 sporadic lymph-node negative patients were
8 preselected, and further filtered by retaining only the genes with a minimum absolute correlation with
9 the disease outcome of 0.3. The resulting list of 231 genes, rank-ordered by absolute correlation,
10 was investigated by sequentially adding the next top 5 genes from the list to a candidate ML classifier
11 and evaluating its performance by leave-one-out cross-validation (LOOCV). This procedure was
12 repeated as long as the estimated accuracy of the classifier improved, providing a final candidate
13 signature of 70 genes. The final signature was validated on multiple independent test sets (including
14 a limited set of 19 external samples in the original study, but several additional validations on
15 independent cohorts in subsequent studies (1,20–23).

16
17
18 The *MammaPrint* signature also provided the role model for the subsequent development of a similar
19 prognostic test for colon cancer, *ColoPrint* (24–29). This test aims at detecting the approx. 20% of
20 patients with stage II colon cancer expected to experience a relapse and develop distant metastases.
21 It uses an 18-gene expression signature, developed by analyzing DNA microarray data in a similar
22 manner to the *MammaPrint* approach. This diagnostic tool has been commercialized as an LDT to
23 assist physicians in selecting treatment options for colon cancer patients. Similar to *MammaPrint*,
24 the development of this signature was characterized by extensive discovery and validation studies,
25 which involved multiple statistical reproducibility, stability and precision analyses across
26 independent, large-scale patient cohorts (30).

27
28
29 Another widely used cancer-related LDT, which received FDA clearance in 2013, is the *Prosigna*
30 *Breast Cancer Prognostic Gene Signature Assay*, previously called *PAM50* test (31–35). This assay
31 assesses mRNA expression for a signature of 58 genes (50 target genes + 8 endogenous control
32 genes) to predict the risk of distant recurrence for hormone-receptor-positive breast cancer from 5
33 to 10 years after diagnosis (prerequisites are that the patients have been treated with hormonal
34 therapy and surgery, and are stage I or stage II lymph-node negative, or in stage II with one to three
35 positive nodes). The development of the test started with a microarray discovery study and involved
36 a multistage filtering approach, using consecutive applications of statistical tests and cross-validation
37 to propose a subset of candidate gene markers (36). The authors compared the reproducibility of
38 classification scores obtained with these markers for three centroid-based prediction methods to
39 ensure the robustness of the methodology. By further developing the approach first into a more
40 sensitive PCR-based test, and then into an assay using the NanoString nCounter Dx Analysis
41 System, the predictive performance was improved in a step-wise fashion. The original discovery
42 study was characterized by significantly larger sample sizes than the majority of reviewed biomarker
43 studies, with a training set of 189 samples, test sets of 761 patients evaluated for prognosis, and
44 133 patients evaluated for prediction of pathologic complete response to treatment with taxane and
45 anthracycline. These study design features in combination with multi-stage filtering and validation
46 approaches, and the improvement of the measurement technology during the course of the study,
47 may explain the successful progression of the PAM50 test to FDA clearance.

48
49
50
51 Among the LDTs for breast cancer prognosis, *Oncotype DX*® is a further test which is already
52 commonly used in clinical practice (3,37–40). The underlying gene signature consists of 16 cancer-
53 associated genes and 5 reference genes, and is therefore often also referred to as '21-gene assay'.
54 Its main application is to predict risk of recurrence in estrogen-receptor positive tumors. The
55 relevance of this prognostic tool for treatment selection is explained by the strong association of the
56 provided recurrence score with the probability of positive treatment response to chemotherapy.
57 *Oncotype DX* was developed using a consecutive refinement procedure, starting with the RT-PCR
58 assessment of 250 candidate genes across 447 patients from three distinct studies to identify the
59 21-gene signature after multiple filtering steps. A recurrence score algorithm built using the signature
60

1
2
3 as input was clinically validated on 668 independent patients (41). The selection of the 16 cancer-
4 related genes included in the assay was mainly done by scoring the performance of the candidate
5 features in all three studies and the consistency of the primer/probe performance in the assay
6 (42). Thus, particular strengths of the development process for this LDT include the consideration of
7 both technical robustness and statistical robustness of the assay across distinct cohorts. However,
8 an independent comparative clinical validation of Oncotype DX and the PAM50 signature for
9 estimating the likelihood of distant recurrence in ER-positive, node-negative, post-menopausal
10 breast cancer patients treated with endocrine therapy suggested that the PAM50 signature provided
11 more predictive information than Oncotype DX (43).
12
13

14 While the first validated omics biomarker signatures were developed for breast cancer, similar
15 diagnostic and prognostic tools have followed for other cancer types. One of these is the Decipher
16 Prostate Cancer Test (4,44–48), which stands out from other omics-derived diagnostic tools in that
17 it is provided together with an additional software platform and database, the Decipher Genomic
18 Resource Information Database (GRID), that captures 1.4 million expression markers per patient to
19 facilitate personalised care. The test itself uses 22 preselected RNAs to predict clinical metastasis
20 and cancer-specific mortality for patients who have undergone radical prostatectomy. An initial
21 discovery study by the Mayo Clinic (Rochester, MN, USA) investigated a cohort of 545 such patients,
22 split into a training (n = 359) and a validation cohort (n = 186). Similar to other LDTs, the discovery
23 started with a genome-wide profiling and used both statistical and ML analyses for filtering. First, *t*-
24 tests were applied (reduction from 1.4 mil. to 18,902 differentially expressed RNAs), then regularized
25 logistic regression (reduction to 43 candidate markers), and finally a random forest-based feature
26 selection (reduction to final set of 22 RNAs). Apart from testing the signature in the validation cohort,
27 further external validations were performed in subsequent studies (4,44–48). Overall, distinctive
28 strengths of the used approach include the improved interpretability of the test results through
29 supporting analyses on the GRID platform, and the robustness of the discovery and validation
30 approach, involving large sample sizes and several complementary statistical and ML assessments.
31
32
33

34 While most diagnostic tests in oncology have been designed for specific cancer types, a dedicated
35 LDT has also been developed for cancers of unknown or uncertain diagnosis. The Cancer Type ID
36 test by bioTheranostics distinguishes between 50 different tumor types using a 92-gene RT-PCR
37 expression measurement signature (10,49–51). This signature has been derived from analyses of a
38 microarray data collection covering 446 frozen tumor samples and 112 formalin-fixed, paraffin-
39 embedded (FFPE) samples of both primary and metastatic tumors. The modeling steps involved *k*-
40 nearest neighbor clustering and classification, and a genetic algorithm to explore the search space
41 of possible feature subset selections. After successful cross-validation (84% accuracy) and external
42 validation (82% accuracy on 112 independent FFPE samples) of the microarray-based signature, it
43 was further developed to use more sensitive RT-PCR measurements. Testing the new approach on
44 an independent validation set provided an increased accuracy (87%). Distinctive characteristics of
45 the overall development process that may have contributed to the positive validation include the
46 efficient and extensive exploration of the search space of possible gene subset selections via a
47 genetic algorithm, the large sample sizes used for discovery and validation, and the transfer of the
48 assay from microarrays to the more sensitive RT-PCR platform.
49
50
51

52 Apart from the omics-derived biomarker signatures that address the most frequent cancer types,
53 more recent applications in oncology focus on the diagnosis of less common malignancies, such as
54 thyroid cancer. Typically, deciding whether a thyroid nodule is benign or cancerous is directly
55 possible via a fine needle aspiration (FNA) biopsy, without requiring more complex measurements
56 or analyses. However, while direct FNA-based diagnosis is feasible in most cases, indeterminate
57 results can occur (52). To help prevent unnecessary surgeries for the corresponding patients, a
58 molecular signature and LDT known as the Afirma™ Gene Expression Classifier (GEC) has been
59 developed to discriminate between benign and cancerous thyroid nodules (52–57). The original
60 discovery study behind the GEC signature used mRNA expression analysis in 315 thyroid nodules,

1
2
3 covering 178 retrospective surgical tissues and 137 prospectively collected FNA samples. The
4 authors trained two ML classifiers separately on surgical tissues and FNAs, assessing the test set
5 performance on 48 independent, prospective FNA samples (50% of which had indeterminate
6 cytopathology). Discriminative features were selected using a linear modeling approach
7 implemented in the software Limma, and a linear support vector machine was applied for model
8 building and performance estimation via 30-fold cross-validation (CV). The successful cross-
9 validation results were later confirmed on multiple distinct cohorts. While the internal validation used
10 in the initial study cannot address cohort-specific biases, the combined use of established feature
11 selection and modeling approaches, and the subsequent external validation across multiple cohorts
12 with large sample sizes may account for the successful translation of this signature.
13
14

15 Most omics-based diagnostic tests identified in the survey rely purely on gene expression profiling
16 data. However, more recently, first multi-omics signatures for diagnostic purposes have been
17 developed. One of the first LDTs that integrated information from both RNA and DNA sequencing
18 was the FoundationOne Heme assay (9,58–60). This assay aims to detect hematologic
19 malignancies, sarcomas, pediatric malignancies, or solid tumors (including among others leukemias,
20 myelodysplastic syndromes, myeloproliferative neoplasms, lymphomas, multiple myeloma, Ewing
21 sarcoma, Leiomyosarcoma, and pediatric tumors). The test identifies four types of genomic
22 alterations (base substitutions, insertions and deletions, copy number alterations, rearrangements)
23 and reports microsatellite instability and tumor mutational burden to facilitate clinical decision
24 making. The approach was originally developed and evaluated using reference samples of pooled
25 cell lines in order to model the main characteristics that determine the test accuracy, including mutant
26 allele frequency, indel length and amplitude of copy change (58). A first validation using 249
27 independent FFPE cancer samples, which had already been characterized by established assays,
28 confirmed the accuracy of the test. Later external validation studies on independent cohorts also
29 corroborated the utility of the test for further diagnostic applications (9,61). The study results highlight
30 the potential of integrating diverse biological data sources in order to obtain more robust and reliable
31 predictions, a strategy that may be promising in particular for complex disorders that involve very
32 heterogeneous phenotypes.
33
34
35
36
37

38 *Non cancer approved omics-derived diagnostic tests (2 Studies)*

39
40 While most clinically approved omics-derived diagnostic tests have been developed in the field of
41 oncology, one of the first LDTs that received FDA clearance for a non-cancer disease was the
42 AlloMap Heart test (8,62–64). It uses a gene expression signature of 11 target genes and 9 control
43 genes in peripheral blood from heart transplant recipients to estimate the risk for acute cellular
44 cardiac allograft rejection. The development process involved statistical analyses of leukocyte
45 microarray profiling data from 285 samples, and subsequent RT-PCR validation and bioinformatics
46 post-processing (8). Prior knowledge from database and literature mining was integrated into the
47 analysis by mapping the data to known alloimmune pathways. This allowed the researchers to
48 narrow down 252 candidate marker genes. An RT-PCR validation on 145 samples confirmed 68 of
49 these candidate genes, which distinguished rejection samples from quiescent samples according to
50 a T-test ($p < 0.01$). Six genes were eliminated due to significant variation in gene expression with
51 sample processing time. Next, correlated gene expression levels were averaged to create robust
52 meta-level features, called 'metagenes', and 20 of these features were added as new variables.
53 Finally, a linear discriminant analysis was applied, providing a prediction model using four individual
54 genes and three metagenes, which aggregate information from 11 original genes. Bootstrap
55 validation procedures and external test set validations were performed to confirm the accuracy of
56 this signature. Overall, distinctive aspects of the development approach for the AlloMap signature
57 include the knowledge-based gene discovery, a comprehensive RT-PCR validation of candidate
58 genes, and the robust bootstrap and external validation analyses.
59
60

1
2
3 The first clinically validated LDT for a cardiovascular indication derived from omics data was the
4 Corus CAD test, developed to identify coronary artery disease (CAD) in stable non-diabetic patients
5 (6,65–68). In contrast to most other omics-based tests, Corus CAD is not a pure molecular signature
6 test, but also takes the clinical covariates gender and age into account. The initial discovery study
7 used a retrospective microarray analysis of blood samples from 195 diabetic and non-diabetic
8 patients from the Duke University CATHGEN registry. After ranking the studied genes in terms of
9 the statistical significance of group differences and prior biological knowledge about their disease
10 relevance, 88 genes were selected for RT-PCR validation. Because diabetes status as a clinical
11 covariate was significantly associated with the observed gene expression alterations, and the
12 identified CAD-associated genes did not overlap between diabetic and non-diabetic patients, the
13 authors decided to limit follow-up work to the non-diabetic patients. In a prospective clinical trial, the
14 PREDICT study, microarray profiling was conducted on blood samples from 198 patients, and top-
15 ranked genes were further validated using RT-PCR for 640 PREDICT blood samples. After multiple
16 filtering steps, taking into account statistical significance in T-tests, biological relevance, gene
17 correlation clustering and cell-type analyses, a final signature of 23 genes was derived, composed
18 of 20 CAD-associated genes and 3 reference genes (69). To maximize the predictive performance,
19 the final prediction algorithm was optimized to adjust for differences associated with age and gender.
20 Overall, compared to other reviewed studies, the Corus CAD approach stands out by taking clinical
21 covariates into account in the final prediction model, including a critical review and adjustment of the
22 inclusion criteria (limiting the focus to nondiabetic patients), and integrating complementary filtering
23 and validation analyses on large sample sizes.
24
25
26
27
28

29 Discussion

31 *Statement of principal findings*

32
33 The scoping review of articles on patient stratification using omics data revealed common limitations
34 in the study design for many published biomarker development projects, such as insufficient and
35 imbalanced sample sizes per study group and inadequate validation methods, but also identified
36 multiple studies that have led to validated diagnostic and prognostic tests. These success stories
37 were investigated in more detail to identify common characteristics in the study design, discovery
38 and validation methods, which may have supported the clinical translation of the initial findings. Key
39 shared aspects that are possible determinants of the study success and could help to guide future
40 biomarker investigations are outlined in Fig. 4. These characteristics, which may serve as a guideline
41 for future studies, cover in particular the following main features:
42
43

44 (1) A sample size selection, study group and replicate design that provides adequate statistical
45 power for the ML analyses;

46
47 2) The application of robust statistical filtering and evaluation schemes (including multiple layers of
48 statistical and ML-based feature selection, combined statistical and biological filters, robust
49 validation schemes that involve multiple cross-validation, bootstrapping and external validation
50 analyses, using multiple suitable and complementary performance metrics, and providing
51 information on the statistical variation and confidence intervals for the performance estimates);
52

53 3) Clarity of the study scope and goals (involving clear inclusion and exclusion criteria, and precisely
54 defined primary and secondary outcomes; new knowledge gained during the project may require a
55 re-definition and adjustment of the inclusion criteria, as in the case of the Corus CAD study, or
56 adjustments in the methodology to reach the goals, such as the progression from non-targeted
57 microarray technology to higher-sensitivity RT-PCR in the case of the Prosigna test and the Cancer
58 Type ID test);
59
60

1
2
3 4) Completeness and reproducibility of the study documentation (covering details on used
4 instruments, parameters and settings, reproducible methods descriptions, and information on data
5 provenance);
6

7 5) Interpretability and biological plausibility of the created predictive models (including explainable
8 and justifiable predictions, human-interpretable model descriptions, and biologically plausible
9 models that agree with the current mechanistic understanding of the studied disorder);
10

11 6) Integration of prior biological knowledge into the predictive feature selection, model building and
12 validation procedures (e.g., using public data on disease-associated molecular pathways and
13 networks; complementary clinical and real-world data, and relevant multi-omics data).
14
15

16 17 *Strengths and Limitations*

18
19 The majority of methodological recommendations derived from the literature survey relate to the
20 early planning and study design for biomarker discovery projects, involving considerations
21 associated with the choice of the study group, sampling and blocking design, the measurement
22 technology, and the input and output variables (11,12). These recommendations are therefore mainly
23 applicable to prospective studies. For retrospective biomarker investigations of already collected
24 data, the suggestions derived from the review are limited to guidance on filtering and evaluation
25 analyses, the integration of prior knowledge from multi-omics data and public annotation databases;
26 and the choice of robust and interpretable modelling approaches for the generation of biologically
27 plausible and reproducible prediction models. While the main focus of the review on studies that
28 have already led to validated biomarker models helps to ensure the quality of the surveyed articles,
29 more recent methodological developments in the ML and cross-validation analysis of omics data,
30 such as meta-learning (70) and bolstered cross-validation (71), will require further dedicated surveys
31 in the future.
32
33
34
35

36 *Discussing important differences in results*

37
38 Previous reviews of ML approaches using omics data for patient stratification have mostly focused
39 on domain-specific analyses for specific types of diseases, or specific types of ML methodologies
40 (72–80). By contrast, this scoping review focused on disease-agnostic workflows with generic
41 applicability across human disorders that involve multifactorial molecular alterations in the affected
42 tissues and body fluids. The coverage of statistical and ML approaches for stratification did not focus
43 on the detailed discussion of specific algorithms or statistics, but rather aimed at identifying key
44 determinants of success for generic analysis workflows that have been applied successfully in
45 practice for biomedical stratification studies.
46
47
48

49 *Meaning of the study: implications for clinicians and policymakers*

50
51 The previous clinical translation successes in omics-based biomarker development reviewed in this
52 survey, which have mostly been achieved for oncology applications, highlight the potential for similar
53 biomarker discovery and validation projects in other fields of biomedicine. In contrast to conventional
54 statistical biomarker discovery approaches, which focus on identifying single-molecule markers,
55 systems-level analysis of omics data using multivariate ML approaches can identify biomarker
56 signatures which are not only more sensitive, specific and robust, but also more biologically insightful
57 in terms of reflecting disease-associated cellular process alterations in a more detailed and
58 comprehensive fashion.
59
60

1
2
3 This scoping review has identified common characteristics of omics studies which have led to
4 clinically validated diagnostic and prognostic tests. Thus, it may help to guide clinical researchers on
5 study design choices, the selection of methodologies for statistical and biological data filtering and
6 ML, and the implementation of adequate validation schemes. By creating awareness of potential
7 pitfalls, such as issues associated with batch effects, biases, confounding factors, lack of statistical
8 power, and multiple hypothesis testing, common reasons for failures in biomarker development can
9 be avoided. This information may also be relevant for policymakers and funding bodies, by facilitating
10 the design of public and private funding schemes for biomedical research in a manner that addresses
11 risks in the funded projects upfront through appropriate guidelines and regulations. Finally, it may
12 help clinicians involved in biomarker discovery to make better use of already available public
13 knowledge and data sources, e.g. cellular pathway and molecular interaction databases, that may
14 allow them to exploit prior knowledge more effectively in biomarker modelling, and create more
15 robust and interpretable prediction models.
16
17
18
19

20 *Unanswered questions & future research*

21
22 Since the recommendations and guidelines identified from the surveyed articles are mostly derived
23 from already established biomarker discovery and validation approaches, new upcoming
24 methodologies could only be covered to a limited extent and may lead to changed recommendations
25 in the future. In particular, in the reviewed patient stratification studies, some of the more recently
26 developed ML and validation workflows, e.g. involving meta-learning, bootstrapping or bolstered
27 cross-validation, are still underrepresented among the reviewed studies, and may play a more
28 important role in the future.
29

30
31 Overall, while the currently available literature on validated stratification biomarkers already provides
32 ample information on common pitfalls and best practices, the development of widely accepted
33 standard guidelines on methodologies for omics biomarker discovery will require further knowledge
34 exchange and deliberation among the stakeholders in the field. In particular, integration of domain-
35 specific expertise in discussions involving clinicians, experimental and data scientists, and regulatory
36 and legal experts is required as a follow-up effort in order to derive comprehensive methodological
37 guidelines for future biomarker development.
38
39
40
41
42

43 **Acknowledgments**

44
45 The authors thank Vanna Pistotti for her assistance with search strategy development and
46 conduction.
47
48

49 **List of Figures**

50
51 Figure 1: Keyword based search strategy for the scoping review

52
53 Figure 2: Study selection flow diagram

54
55 Figure 3: Validation methods used in omics biomarker studies

56
57 Figure 4: Characteristics of successful omics-based studies
58
59
60

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Definitions (In boxes)

Box 1: *What is Personalised Medicine?*

According to the European Council Conclusion on personalised medicine for patients, personalised medicine is 'a medical model using characterisation of individuals' phenotypes and genotypes (e.g., molecular profiling, medical imaging, lifestyle data) for tailoring the right therapeutic strategy for the right person at the right time, and/or to determine the predisposition to disease and/or to deliver timely and targeted prevention (81).

In the context of the Permit project, we applied the following common operational definition of personalised medicine research: a set of comprehensive methods (methodology, statistics, validation, technology) to be applied in the different phases of the development of a personalised approach to treatment, diagnosis, prognosis, or risk prediction. Ideally, robust and reproducible methods should cover all the steps between the generation of the hypothesis (e.g., a given stratum of patients could better respond to a treatment), its validation and pre-clinical development, and up to the definition of its value in a clinical setting (16).

References

1. Cardoso F, van't Veer LJ, Bogaerts J, Slaets L, Viale G, Delaloge S, et al. 70-Gene Signature as an Aid to Treatment Decisions in Early-Stage Breast Cancer. *N Engl J Med*. 2016;375(8):717–29.
2. Bachleitner-Hofmann T, Simon I, Salazar R, Tabernero J, Rosenberg R, van der Akker J, et al. Development and Validation of a Robust Molecular Diagnostic Test (COLOPRINT) for Predicting Outcome in Stage II Colon Cancer Patients. *Ann Oncol*. 2012;
3. Knezevic D, Goddard AD, Natraj N, Cherbavaz DB, Clark-Langone KM, Snable J, et al. Analytical validation of the Oncotype DX prostate cancer assay - a clinical RT-PCR assay optimized for prostate needle biopsies. *BMC Genomics*. 2013;14(1).
4. Torres A, Alshalalfa M, Tomlins SA, Erho N, Gibb EA, Chelliserry J, et al. Comprehensive Determination of Prostate Tumor ETS Gene Status in Clinical Samples Using the CLIA Decipher Assay. *J Mol Diagnostics*. 2017;19(3):475–84.
5. Angell TE, Babiarz J, Barth N, Blevins T, Duh Q, Ghossein RA, et al. Clinical validation of the AFIRMA genomic sequencing braf V600E classifier. *Thyroid [Internet]*. 2017;27:A50. Available from: <http://www.embase.com/search/results?subaction=viewrecord&from=export&id=L624116485>
6. Ladapo JA, Budoff MJ, Sharp D, Kuo JZ, Huang L, Maniet B, et al. Utility of a Precision Medicine Test in Elderly Adults with Symptoms Suggestive of Coronary Artery Disease. *J Am Geriatr Soc*. 2018;66(2):309–15.
7. Tabari E, Lovejoy AF, Lin H, Bolen CR, Saelee SL, Lefkowitz JP, et al. Molecular characteristics and disease burden metrics determined by next-generation sequencing on circulating tumor DNA correlate with progression free survival in previously untreated diffuse large B-cell lymphoma. *Blood [Internet]*. 2019;134. Available from: <http://dx.doi.org/10.1182/blood-2019-123633>
8. Deng MC. The AlloMap™ genomic biomarker story: 10 years after. *Clin Transplant*. 2017;31(3).
9. He J, Abdel-Wahab O, Nahas MK, Wang K, Rampal RK, Intlekofer AM, et al. Integrated genomic DNA/RNA profiling of hematologic malignancies in the clinical setting. *Blood*. 2016;127(24):3004–14.
10. Ma XJ, Patel R, Wang X, Salunga R, Murage J, Desai R, et al. Molecular classification of human cancers using a 92-gene real-time quantitative polymerase chain reaction assay. *Arch Pathol Lab Med*. 2006;130(4):465–73.
11. McShane LM, Cavenagh MM, Lively TG, Eberhard DA, Bigbee WL, Williams PM, et al. Criteria for the use of omics-based predictors in clinical trials. *Nature*. 2013.
12. Moons KGM, Altman DG, Reitsma JB, Ioannidis JPA, Macaskill P, Steyerberg EW, et al. Transparent reporting of a multivariable prediction model for individual prognosis or diagnosis (TRIPOD): Explanation and elaboration. *Ann Intern Med*. 2015;162(1):W1–73.
13. Vinet L, Zhedanov A. A “missing” family of classical orthogonal polynomials. *J Phys A Math Theor*. 2011;44(8).
14. Pham MT, Rajić A, Greig JD, Sargeant JM, Papadopoulos A, Mcewen SA. A scoping review of scoping reviews: Advancing the approach and enhancing the consistency. *Res Synth Methods*. 2014;5(4):371–85.
15. Munn Z, Peters MDJ, Stern C, Tufanaru C, McArthur A, Aromataris E. Systematic review or scoping review? Guidance for authors when choosing between a systematic or scoping

- review approach. *BMC Med Res Methodol*. 2018;18(1):143.
16. Banzi R, Gerardi C, Fratelli M, Garcia P, Torres T, Abad JMH, et al. Methodological approaches for personalised medicine: protocol for a series of scoping reviews [Internet]. 10.5281/zenodo.3770937. Available from: <https://zenodo.org/record/3770937>
 17. Tricco AC, Lillie E, Zarin W, O'Brien KK, Colquhoun H, Levac D, et al. PRISMA extension for scoping reviews (PRISMA-ScR): Checklist and explanation. *Ann Intern Med*. 2018 Oct;169(7):467–73.
 18. Ouzzani M, Hammady H, Fedorowicz Z, Elmagarmid A. Rayyan-a web and mobile app for systematic reviews. *Syst Rev*. 2016;5(1):210.
 19. Van't Veer LJ, Dai H, Van de Vijver MJ, He YD, Hart AAM, Mao M, et al. Gene expression profiling predicts clinical outcome of breast cancer. *Nature*. 2002;415(6871):530–6.
 20. Winner BS, Sgroi DC, Ryan PD, Bruinsma TJ, Glas AM, Male A, et al. Analysis of the mamma print breast cancer assay in a predominantly postmenopausal cohort. *Clin Cancer Res*. 2008;14(10):2988–93.
 21. Slodkowska EA, Ross JS. MammaPrint™ 70-gene signature: Another milestone in personalized medical care for breast cancer patients. *Expert Rev Mol Diagn*. 2009;9(5):417–22.
 22. Sapino A, Roepman P, Linn SC, Snel MHJ, Delahaye LJM, Van Den Akker J, et al. MammaPrint molecular diagnostics on formalin-fixed, paraffin-embedded tissue. *J Mol Diagnostics*. 2014;16(2):190–7.
 23. Mook S, Knauer M, Bueno-De-Mesquita JM, Retel VP, Wesseling J, Linn SC, et al. Metastatic potential of T1 breast cancer can be predicted by the 70-gene MammaPrint signature. *Ann Surg Oncol*. 2010;17(5):1406–13.
 24. Maak M, Simon I, Nitsche U, Roepman P, Snel M, Glas AM, et al. Independent validation of a prognostic genomic signature (ColoPrint) for patients with stage II colon cancer. *Ann Surg*. 2013;257(6):1053–8.
 25. Kopetz S, Taberero J, Rosenberg R, Jiang Z, Moreno V, Bachleitner-Hofmann T, et al. Genomic Classifier ColoPrint Predicts Recurrence in Stage II Colorectal Cancer Patients More Accurately Than Clinical Factors. *Oncologist*. 2015;20(2):127–33.
 26. Tan IB, Tan P. Genetics: An 18-gene signature (ColoPrint®) for colon cancer prognosis. *Nat Rev Clin Oncol*. 2011;8(3):131–3.
 27. Rosenberg R, Maak M, Simon I, Nitsche U, Schuster T, Kuenzli B, et al. Independent validation of a prognostic genomic profile (ColoPrint) for stage II colon cancer (CC) patients. *J Clin Oncol*. 2011;29(4_suppl):358–358.
 28. Salazar R, de Waard JW, Glimelius B, Marshall J, Klaase J, Van Der Hoeven J, et al. The PARSC trial, a prospective study for the assessment of recurrence risk in stage II colon cancer (CC) patients using ColoPrint. *J Clin Oncol*. 2012;30(4_suppl):678–678.
 29. Taberero J, Moreno V, Rosenberg R, Nitsche U, Bachleitner-Hofmann T, Lanza G, et al. Clinical and technical validation of a genomic classifier (ColoPrint) for predicting outcome of patients with stage II colon cancer. *J Clin Oncol*. 2012;30(4_suppl):384–384.
 30. Bachleitner-Hofmann T, Simon I, Salazar R, Taberero J, Rosenberg R, van der Akker J, et al. Development and Validation of a Robust Molecular Diagnostic Test (COLOPRINT) for Predicting Outcome in Stage II Colon Cancer Patients. *Ann Oncol*. 2012;23:ix179.
 31. Nielsen T, Wallden B, Schaper C, Ferree S, Liu S, Gao D, et al. Analytical validation of the PAM50-based Prosigna Breast Cancer Prognostic Gene Signature Assay and nCounter Analysis System using formalin-fixed paraffin-embedded breast tumor specimens. *BMC*

- Cancer. 2014;14(1).
32. Wallden B, Storhoff J, Nielsen T, Dowidar N, Schaper C, Ferree S, et al. Development and verification of the PAM50-based Prosigna breast cancer gene signature assay. *BMC Med Genomics*. 2015;8(1).
 33. Alvarado MD, Prasad C, Rothney M, Cherbavaz DB, Sing AP, Baehner FL, et al. A Prospective Comparison of the 21-Gene Recurrence Score and the PAM50-Based Prosigna in Estrogen Receptor-Positive Early-Stage Breast Cancer. *Adv Ther*. 2015;32(12):1237–47.
 34. Jensen MB, Lænkholm AV, Nielsen TO, Eriksen JO, Wehn P, Hood T, et al. The Prosigna gene expression assay and responsiveness to adjuvant cyclophosphamide-based chemotherapy in premenopausal high-risk patients with breast cancer. *Breast Cancer Res*. 2018;20(1).
 35. Hequet D, Callens C, Gentien D, Albaud B, Mouret-Reynier MA, Dubot C, et al. Prospective, multicenter French study evaluating the clinical impact of the Breast Cancer Intrinsic Subtype-Prosigna® Test in the management of early-stage breast cancers. *PLoS One*. 2017;12(10).
 36. Parker JS, Mullins M, Cheang MCU, Leung S, Voduc D, Vickery T, et al. Supervised risk predictor of breast cancer based on intrinsic subtypes. *J Clin Oncol*. 2009;27(8):1160–7.
 37. Kelly CM, Krishnamurthy S, Bianchini G, Litton JK, Gonzalez-Angulo AM, Hortobagyi GN, et al. Utility of oncotype DX risk estimates in clinically intermediate risk hormone receptor-positive, HER2-normal, grade II, lymph node-negative breast cancers. *Cancer*. 2010;116(22):5161–7.
 38. Lo SS, Mumby PB, Norton J, Rychlik K, Smerage J, Kash J, et al. Prospective multicenter study of the impact of the 21-gene recurrence score assay on medical oncologist and patient adjuvant breast cancer treatment selection. *J Clin Oncol*. 2010;28(10):1671–6.
 39. Carlson JJ, Roth JA. The impact of the Oncotype Dx breast cancer assay in clinical practice: A systematic review and meta-analysis. Vol. 141, *Breast Cancer Research and Treatment*. 2013. p. 13–22.
 40. Thakur SS, Li H, Chan AMY, Tudor R, Bigras G, Morris D, et al. The use of automated Ki67 analysis to predict Oncotype DX risk-of-recurrence categories in early-stage breast cancer. *PLoS One*. 2018/01/06. 2018;13(1):e0188983.
 41. Gianni L, Zambetti M, Clark K, Baker J, Cronin M, Wu J, et al. Gene expression profiles in paraffin-embedded core biopsy tissue predict response to chemotherapy in women with locally advanced breast cancer. *J Clin Oncol*. 2005;23(29):7265–77.
 42. Paik S, Shak S, Tang G, Kim C, Baker J, Cronin M, et al. A multigene assay to predict recurrence of tamoxifen-treated, node-negative breast cancer. *N Engl J Med*. 2004;351(27):2817–26.
 43. Dowsett M, Sestak I, Lopez-Knowles E, Sidhu K, Dunbier AK, Cowens JW, et al. Comparison of PAM50 risk of recurrence score with oncotype DX and IHC4 for predicting risk of distant recurrence after endocrine therapy. *J Clin Oncol*. 2013;31(22):2783–90.
 44. Marrone M, Potosky AL, Penson D, Freedman AN. A 22 gene-expression assay, decipher® (GenomeDx biosciences) to predict five-year risk of metastatic prostate cancer in men treated with radical prostatectomy. *PLoS Curr*. 2015;7(EVIDENCEONGENOMICTESTS).
 45. Nguyen PL, Haddad Z, Lam LLC, Ong K, Buerki C, Deheshi S, et al. Evaluation of the Decipher prostate cancer classifier to predict metastasis and disease-specific mortality from genomic analysis of diagnostic prostate needle biopsy specimens. *J Clin Oncol*. 2017;35(6_suppl):4–4.

- 1
- 2
- 3 46. Magi-Galluzzi C, Yousefi K, Haddad Z, Palmer-Aronsten B, Lam LLC, Buerki C, et al.
- 4 Validation of the Decipher prostate cancer classifier for predicting 10-year postoperative
- 5 metastasis from analysis of diagnostic needle biopsy specimens. *J Clin Oncol*.
- 6 2016;34(2_suppl):59–59.
- 7
- 8 47. Dalela D, Löppenber B, Sood A, Sammon J, Abdollah F. Contemporary Role of the
- 9 Decipher® Test in Prostate Cancer Management: Current Practice and Future Perspectives.
- 10 *Rev Urol*. 2016;18(1):1–9.
- 11
- 12 48. Klein EA, Haddad Z, Yousefi K, Lam LLC, Wang Q, Choeurng V, et al. Decipher Genomic
- 13 Classifier Measured on Prostate Biopsy Predicts Metastasis Risk. *Urology*. 2016;90:148–52.
- 14
- 15 49. Weiss LM, Chu P, Schroeder BE, Singh V, Zhang Y, Erlander MG, et al. Blinded comparator
- 16 study of immunohistochemical analysis versus a 92-gene cancer classifier in the diagnosis
- 17 of the primary site in metastatic tumors. *J Mol Diagnostics*. 2013;15(2):263–9.
- 18
- 19 50. Greco FA, Spigel DR, Yardley DA, Erlander MG, Ma X, Hainsworth JD. Molecular Profiling
- 20 in Unknown Primary Cancer: Accuracy of Tissue of Origin Prediction. *Oncologist*.
- 21 2010;15(5):500–6.
- 22
- 23 51. Hainsworth JD, Rubin MS, Spigel DR, Boccia R V., Raby S, Quinn R, et al. Molecular gene
- 24 expression profiling to predict the tissue of origin and direct site-specific therapy in patients
- 25 with carcinoma of unknown primary site: A prospective trial of the Sarah cannon research
- 26 institute. *J Clin Oncol*. 2013;31(2):217–23.
- 27
- 28 52. Harrison G, Sosa JA, Jiang X. Evaluation of the afirma gene expression classifier in repeat
- 29 indeterminate thyroid nodules. *Arch Pathol Lab Med*. 2017;141(7):985–9.
- 30
- 31 53. Chudova D, Wilde JI, Wang ET, Wang H, Rabbee N, Egidio CM, et al. Molecular
- 32 classification of thyroid nodules using high-dimensionality genomic data. *J Clin Endocrinol*
- 33 *Metab*. 2010;95(12):5296–304.
- 34
- 35 54. Kim MI, Alexander EK. Diagnostic use of molecular markers in the evaluation of thyroid
- 36 nodules. Vol. 18, *Endocrine Practice*. 2012. p. 796–802.
- 37
- 38 55. Ali SZ, Fish SA, Lanman R, Randolph GW, Sosa JA. Use of the Afirma® gene expression
- 39 classifier for preoperative identification of benign thyroid nodules with indeterminate fine
- 40 needle aspiration cytopathology. *PLoS Currents*. 2013. p. 1–7.
- 41
- 42 56. McIver B, Castro MR, Morris JC, Bernet V, Smallridge R, Henry M, et al. An independent
- 43 study of a gene expression classifier (Afirma) in the evaluation of cytologically indeterminate
- 44 thyroid nodules. *J Clin Endocrinol Metab*. 2014;99(11):4069–77.
- 45
- 46 57. Lastra RR, Pramick MR, Crammer CJ, LiVolsi VA, Baloch ZW. Implications of a suspicious
- 47 afirma test result in thyroid fine-needle aspiration cytology: An institutional experience.
- 48 *Cancer Cytopathol*. 2014;122(10):737–44.
- 49
- 50 58. Frampton GM, Fichtenholtz A, Otto GA, Wang K, Downing SR, He J, et al. Development
- 51 and validation of a clinical cancer genomic profiling test based on massively parallel DNA
- 52 sequencing. *Nat Biotechnol*. 2013;31(11):1023–31.
- 53
- 54 59. Wang K, Sanchez-Martin M, Wang X, Knapp KM, Koche R, Vu L, et al. Patient-derived
- 55 xenotransplants can recapitulate the genetic driver landscape of acute leukemias.
- 56 *Leukemia*. 2017;31(1):151–8.
- 57
- 58 60. Tarlock K, He J, Zhong S, Ries RE, Bailey M, Morley S, et al. Distinct age-associated
- 59 genomic profiles in acute myeloid leukemia (AML) using FoundationOne heme. *J Clin*
- 60 *Oncol*. 2016;34(15_suppl):7041–7041.
61. Lieber DS, Kennedy MR, Johnson DB, Greenbowe JR, Frampton GM, Schrock AB, et al.
- Abstract B16: Validation and clinical feasibility of a Foundation Medicine assay to identify

- immunotherapy response potential through tumor mutational burden (TMB). In 2017.
62. Yamani MH, Taylor DO, Rodriguez ER, Cook DJ, Zhou L, Smedira N, et al. Transplant Vasculopathy Is Associated With Increased AlloMap Gene Expression Score. *J Hear Lung Transplant*. 2007;26(4):403–6.
 63. Yamani MH, Taylor DO, Haire C, Smedira N, Starling RC. Post-transplant ischemic injury is associated with up-regulated AlloMap gene expression. *Clin Transplant*. 2007;21(4):523–5.
 64. Kobashigawa J, Patel J, Azarbal B, Kittleson M, Chang D, Czer L, et al. Randomized Pilot Trial of Gene Expression Profiling Versus Heart Biopsy in the First Year after Heart Transplant: Early Invasive Monitoring Attenuation Through Gene Expression Trial. *Circ Hear Fail*. 2015;8(3):557–64.
 65. Wingrove JA, Daniels SE, Sehnert AJ, Tingley W, Elashoff MR, Rosenberg S, et al. Correlation of peripheral-blood gene expression with the extent of coronary artery stenosis. *Circ Cardiovasc Genet*. 2008;1(1):31–8.
 66. Rosenberg S, Dehais C, Ducray F, Alentron A, Tanguy M, De Reyneis A, et al. OS11.3 Machine learning for better prognostic stratification and driver genes identification in 1p/19q-codeleted grade III gliomas. *Neuro Oncol* [Internet]. 2017;19(suppl_3):iii22–iii22. Available from: <http://dx.doi.org/10.1093/neuonc/nox036>
 67. Vargas J, Lima JAC, Kraus WE, Douglas PS, Rosenberg S. Use of the Corus® CAD Gene Expression Test for Assessment of Obstructive Coronary Artery Disease Likelihood in Symptomatic Non-Diabetic Patients. *PLoS Currents*. 2013.
 68. Elashoff MR, Wingrove JA, Beineke P, Daniels SE, Tingley WG, Rosenberg S, et al. Development of a blood-based gene expression algorithm for assessment of obstructive coronary artery disease in non-diabetic patients. *BMC Med Genomics*. 2011;4.
 69. Rosenberg S, Elashoff MR, Beineke P, Daniels SE, Wingrove JA, Tingley WG, et al. Multicenter validation of the diagnostic accuracy of a blood-based gene expression test for assessing obstructive coronary artery disease in nondiabetic patients. *Ann Intern Med*. 2010;153(7):425–34.
 70. Lemke C, Budka M, Gabrys B. Metalearning: a survey of trends and technologies. *Artif Intell Rev*. 2015;
 71. Sima C, Braga-Neto UM, Dougherty ER. High-dimensional bolstered error estimation. *Bioinformatics*. 2011;
 72. Ali M, Aittokallio T. Machine learning and feature selection for drug response prediction in precision oncology applications. *Biophysical Reviews*. 2019.
 73. Barla A, Jurman G, Riccadonna S, Merler S, Chierici M, Furlanello C. Machine learning methods for predictive proteomics. *Br Bioinform*. 2008/03/04. 2008;9(2):119–28.
 74. Grollemund V, Pradat PF, Querin G, Delbot F, Le Chat G, Pradat-Peyre JF, et al. Machine learning in amyotrophic lateral sclerosis: Achievements, pitfalls, and future directions. *Front Neurosci* [Internet]. 2019;13. Available from: <http://dx.doi.org/10.3389/fnins.2019.00135>
 75. Ho DSW, Schierding W, Wake M, Saffery R, O'Sullivan J. Machine learning SNP based prediction for precision medicine. *Front Genet* [Internet]. 2019;10(MAR). Available from: <http://dx.doi.org/10.3389/fgene.2019.00267>
 76. Krittanawong C, Bomback AS, Baber U, Bangalore S, Messerli FH, Wilson Tang WH. Future Direction for Using Artificial Intelligence to Predict and Manage Hypertension. *Curr Hypertens Rep*. 2018/07/08. 2018;20(9):75.
 77. Long NP, Yoon SJ, Anh NH, Nghi TD, Lim DK, Hong YJ, et al. A systematic review on metabolomics-based diagnostic biomarker discovery and validation in pancreatic cancer.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Metabolomics. 2019/03/05. 2018;14(8):109.

78. Martinez BI, Stabenfeldt SE. Current trends in biomarker discovery and analysis tools for traumatic brain injury. *J Biol Eng* [Internet]. 2019;13(1). Available from: <http://dx.doi.org/10.1186/s13036-019-0145-8>
79. Patil S, Awan KH, Arakeri G, Seneviratne CJ, Muddur N, Malik S, et al. Machine learning and its potential applications to the genomic study of head and neck cancer-A systematic review. *J Oral Pathol Med*. 2019;48(9):773–9.
80. Saini G, Mittal K, Rida P, Janssen EAM, Gogineni K, Aneja R. Panoptic view of prognostic models for personalized breast cancer management. *Cancers (Basel)* [Internet]. 2019;11(9). Available from: <http://dx.doi.org/10.3390/cancers11091325>
81. European Council. Council conclusions on personalised medicine for patients. *Off J Eur Union* [Internet]. 2015;431(2):1–4. Available from: [https://eur-lex.europa.eu/legal-content/EN/TXT/PDF/?uri=CELEX:52015XG1217\(01\)&from=EN](https://eur-lex.europa.eu/legal-content/EN/TXT/PDF/?uri=CELEX:52015XG1217(01)&from=EN)

Author Contributions

Study conception and design: EG, AR.

Methodology: CG, RB.

Data collection and analysis: EG, AR

Original draft preparation: EG

Review and editing: AR, PG, CG, JDM, RB.

Project supervision: PG

Funding acquisition: JDM.

All authors have read and revised the manuscript and approved the final version.

The members of the PERMIT group were involved in the preparation or revision of the joint protocol of the four scoping reviews of the PERMIT series, attended the joint workshop (consultation exercise) and are co-authors of the other scoping reviews of the PERMIT series.

Collaborators

PERMIT group:

1. Antonio L. Andreu
2. Florence Bietrix,
3. Florie Brion Bouvier
4. Montserrat Carmona Rodriguez
5. Maria del Mar Polo-de Santos,
6. Maddalena Fratelli,
7. Rainer Girgenrath,
8. Alexander Grundmann,
9. Josep Maria Haro,
10. Frank Hulstaert,
11. Iñaki Imaz-Iglesia,
12. Setefilla Luengo Matos
13. Emmet McCormack,
14. Albert Sanchez Niubo,
15. Emanuela Oldoni,
16. Raphael Porcher,
17. Vibeke Fosse,
18. Luis M. Sánchez-Gómez,
19. Lorena San Miguel,
20. Cecilia Superchi,
21. Teresa Torres,
22. Anna Monistrol Mula

Funding statement

This project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No. 874825.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Competing interests statement

None declared

Ethics approval

This study was based entirely on a systematic review of relevant published literature and did not require an ethics approval.

Patient consent

This study did not require consent from patients as it uses no individual data.

Permission to reproduce material from other sources

This study has cited all references which are published and publicly available.

Data sharing statement

The study protocol was published on the online platform Zenodo (<https://zenodo.org/record/3770937>). Copies of searches and data extraction sheets will be made publicly available on Zenodo as part of the database collection for all scoping reviews conducted in the PERMIT project.

Figure legends

Fig. 1. Keyword based search strategy for the scoping review. Four categories of keywords were defined to retrieve relevant articles from the biomedical literature on machine learning analyses of omics data for personalised medicine, which include a validation study (highlighted by the colored boxes in the center). For each category relevant keywords were determined, including controlled vocabulary terms from the Medical Subject Headings (MeSH) thesaurus by the US National Library of Medicine (upper and lower boxes with frames colored according to the corresponding category). As indicated by the keyword “AND” in the center, a conjunctive search was conducted, i.e., every retrieved article had to contain at least one keyword from each category. This strategy was adapted for searching the other databases.

Fig. 2. Study selection flow diagram. Flow diagram of the procedure for the scoping review article identification, screening, eligibility assessment, and final inclusion, according to the PRISMA scheme (17). Reasons for excluding full-text were not mutually exclusive.

Fig. 3. Validation methods used in omics biomarker studies. Stacked bar chart of the number of articles retrieved in the scoping review for different categories of validation methods used in the underlying biomarker studies (covering time periods from 2000 to 2020). The majority of studies use only internal cohort validation approaches, such as cross-validation (CV), training/test set split validation, out-of-bag validation (for tree-based classifiers), and combinations of CV and test set validation within the same cohort. Studies with an external validation on an independent patient cohort (with or without an additional internal cross-validation) are still underrepresented, even in more recent time periods. All filtered full-text articles derived from the scoping review except for review articles were included in the analysis.

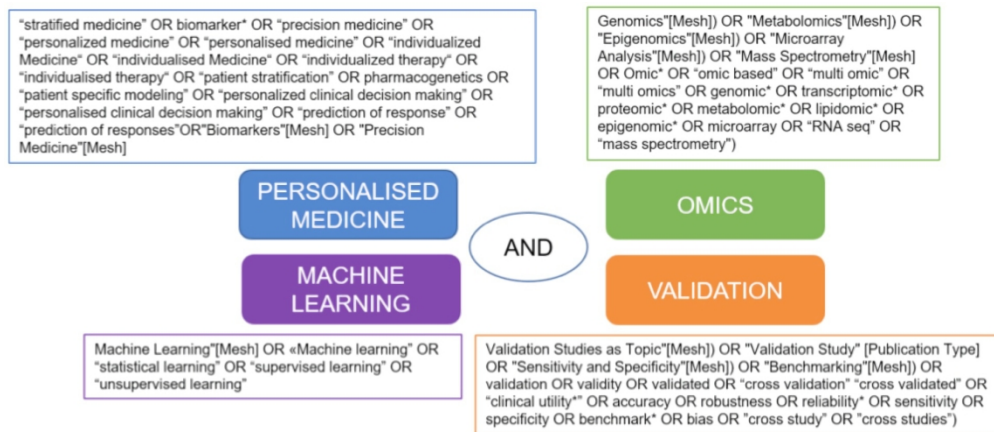
Fig. 4. Characteristics of successful omics-based studies. Six main categories of design and implementation aspects that characterize successful omics-based biomarker development studies were identified (starting from the centre left in the figure and proceeding clockwise): 1) Adequacy of the study design & sample size selection; 2) Rigor and robustness of the statistical evaluation; 3) Clarity of scope and goals; 4) Completeness and reproducibility of the study documentation; 5) Interpretability and biological plausibility of the created predictive models; 6) Integration of prior biological knowledge into the model building and validation procedures.

1
2
3 **Tables**
4
5
6

7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41

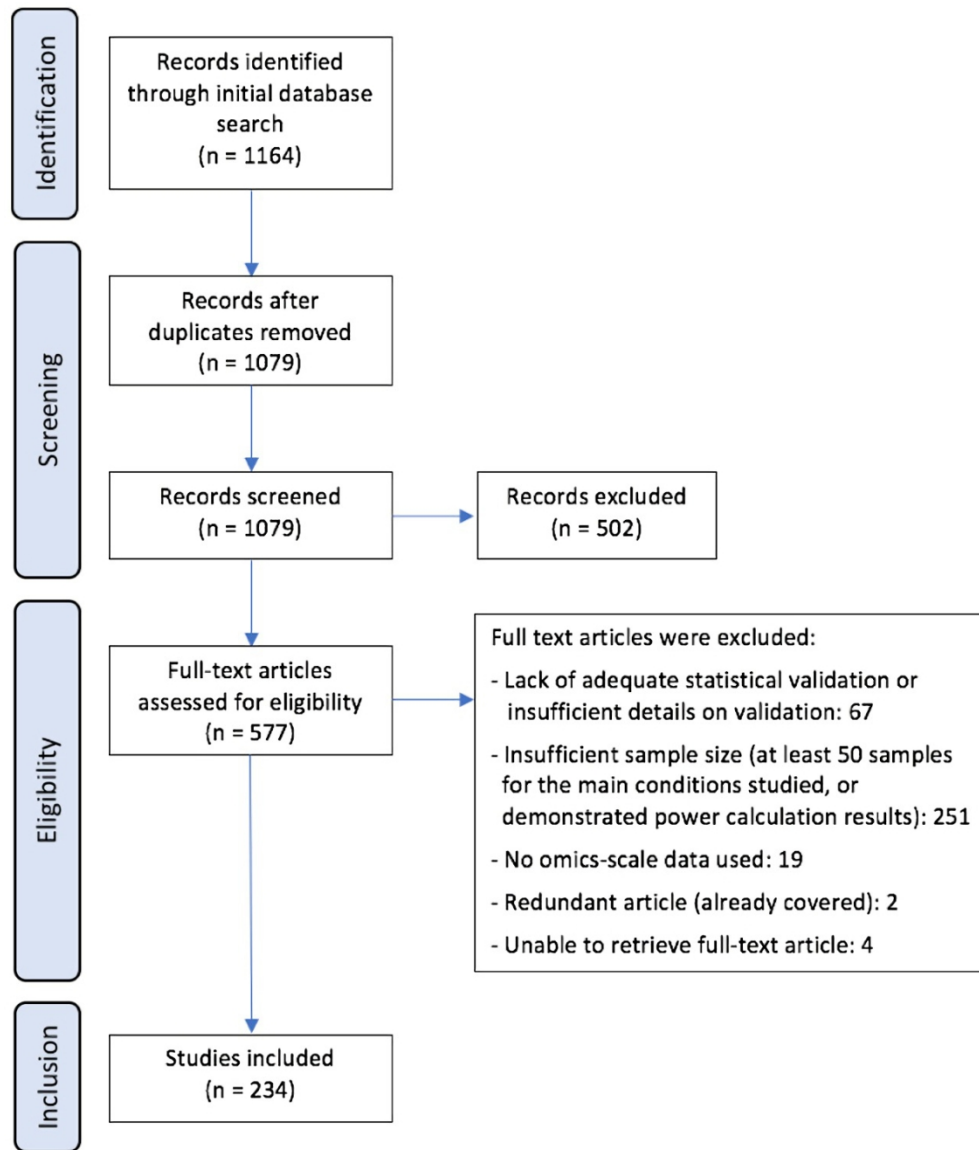
Name	Test approval (FDA-cleared and/or LDT)	Purpose	References
MammaPrint	FDA-cleared, LDT	breast cancer risk-of-recurrence assessment	(1,20–23)
ColoPrint	LDT	colon cancer development of distant metastasis prediction	(24–29)
Prosigna Assay / PAM50	FDA-cleared, LDT	breast cancer risk of distant recurrence prediction	(31–35)
Oncotype DX	LDT	breast cancer risk-of-recurrence assessment	(3,37–40)
Decipher	LDT	prostate cancer metastatic risk prediction	(4,44–48)
Cancer Type ID	LDT	predict tumor type for cancers of unknown / uncertain diagnosis	(10,49–51)
Afirma™ Gene Expression Classifier	LDT	discriminate between benign and cancerous thyroid nodules	(52–57)
Foundation One Heme	LDT	test for hematologic malignancies, sarcomas, or solid tumors	(9,58–60)
AlloMap Heart	FDA-cleared, LDT	identifying heart transplant recipients with risk of cellular rejection	(8,62–64)
Corus CAD	LDT	identify obstructive coronary artery disease	(6,65–68)

42
43 **Tab. 1.** Examples of clinically approved omics-derived diagnostic or prognostic tests designs applied
44 to personalised medicine (synonyms for the same test are separated by the “/”-symbol). FDA-
45 approval status was checked on the web-site [https://www.fda.gov/medical-devices/vitro-](https://www.fda.gov/medical-devices/vitro-diagnostics/nucleic-acid-based-tests)
46 [diagnostics/nucleic-acid-based-tests](https://www.fda.gov/medical-devices/vitro-diagnostics/nucleic-acid-based-tests) and reflects the status as of Oct. 22nd 2020.
47
48
49
50
51
52
53
54
55
56
57
58
59
60



Keyword based search strategy for the scoping review. Four categories of keywords were defined to retrieve relevant articles from the biomedical literature on machine learning analyses of omics data for personalised medicine, which include a validation study (highlighted by the colored boxes in the center). For each category relevant keywords were determined, including controlled vocabulary terms from the Medical Subject Headings (MeSH) thesaurus by the US National Library of Medicine (upper and lower boxes with frames colored according to the corresponding category). As indicated by the keyword "AND" in the center, a conjunctive search was conducted, i.e., every retrieved article had to contain at least one keyword from each category. This strategy was adapted for searching the other databases.

164x71mm (300 x 300 DPI)

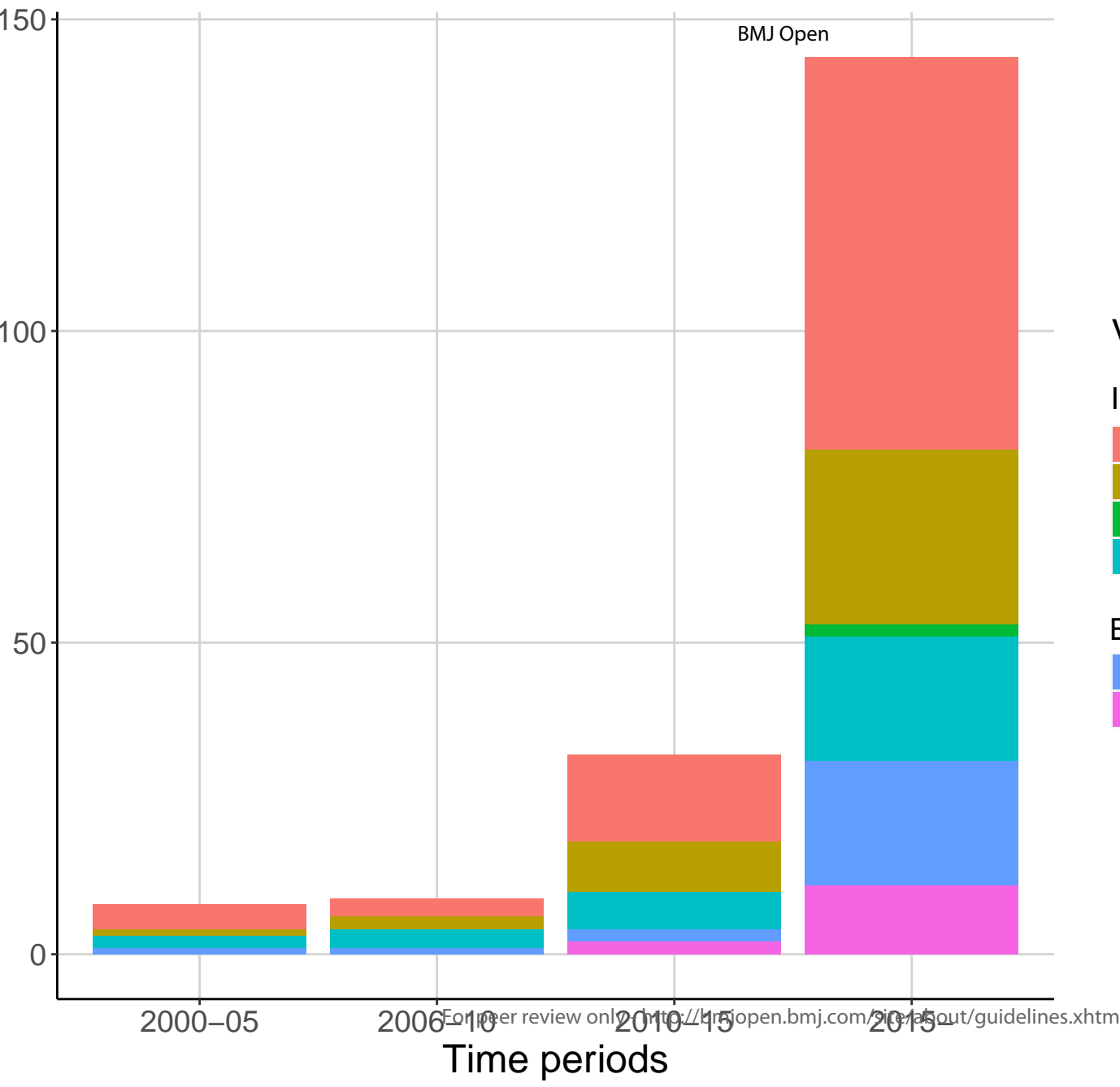


43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Study selection flow diagram. Flow diagram of the procedure for the scoping review article identification, screening, eligibility assessment, and final inclusion, according to the PRISMA scheme. Reasons for excluding full-text were not mutually exclusive.

109x127mm (300 x 300 DPI)

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38



Validation methods

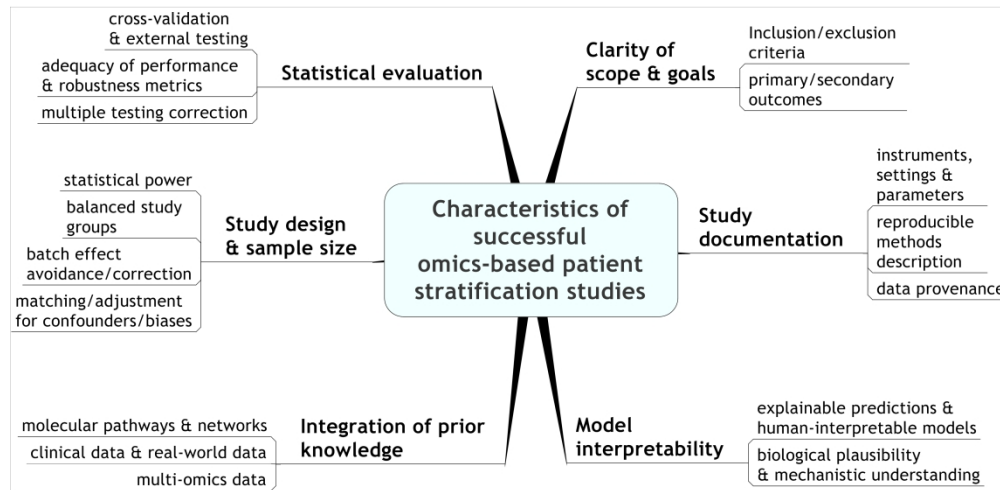
Internal validation:

- Cross-validation (CV)
- Training/test set validation
- Out-of-bag validation
- CV + internal cohort validation

External validation:

- CV + external cohort validation
- External cohort validation

on 6 December 2021. Downloaded from <http://bmjopen.bmj.com/> on April 8, 2022 by guest. Protected by copyright.



Characteristics of successful omics-based studies. Six main categories of design and implementation aspects that characterize successful omics-based biomarker development studies were identified (starting from the centre left in the figure and proceeding clockwise): 1) Adequacy of the study design & sample size selection; 2) Rigor and robustness of the statistical evaluation; 3) Clarity of scope and goals; 4) Completeness and reproducibility of the study documentation; 5) Interpretability and biological plausibility of the created predictive models; 6) Integration of prior biological knowledge into the model building and validation procedures.

338x165mm (300 x 300 DPI)

Online Supplementary file 2 – Search strategy

Precise queries and number of items retrieved for each query for the keyword searches conducted in the databases *PubMed*, *EMBASE* and *Web of Science* as part of the scoping review.

1) PubMed Query

No.	Query	Items found
#14	Search #4 AND #7 AND #10 AND #13 Sort by: PublicationDate	365
#13	Search #11 OR #12 Sort by: PublicationDate	2480122
#12	Search (Validation Studies as Topic [Mesh]) OR "Validation Study" [Publication Type] OR "Sensitivity and Specificity" [Mesh]) OR "Benchmarking" [Mesh])	0
#11	Search (validation OR validity OR validated OR "cross validation" "cross validated" OR "clinical utility*" OR accuracy OR robustness OR reliability* OR sensitivity OR specificity OR benchmark* OR bias OR "cross study" OR "cross studies")	2480122
#10	Search #8 OR #9 Sort by: PublicationDate	982942
#9	Search ("Genomics" [Mesh]) OR "Metabolomics" [Mesh]) OR "Epigenomics" [Mesh]) OR "Microarray Analysis" [Mesh]) OR "Mass Spectrometry" [Mesh])	422277
#8	Search (Omic* OR "omic based" OR "multi omic" OR "multi omics" OR genomic* OR transcriptomic* OR proteomic* OR metabolomic* OR lipidomic* OR epigenomic* OR microarray OR "RNA seq" OR "mass spectrometry")	944832
#7	Search #5 OR #6 Sort by: PublicationDate	743829
#6	Search ("Biomarkers" [Mesh]) OR "Precision Medicine" [Mesh])	743829
#5	Search (stratified medicine" OR cluster* OR "sub group*" OR Subgroup* OR biomarker* OR diagnos* OR prognos* OR "precision medicine" OR "personalized medicine"OR "personalised medicine" OR "individualized Medicine" OR "individualised Medicine" OR "individualized therapy" OR "individualised therapy")	0
#4	Search #2 OR #3 Sort by: PublicationDate	40640
#3	Search "Machine Learning" [Mesh] Sort by: PublicationDate	16481
#2	Search ("Machine learning" OR "statistical learning" OR "supervised learning" OR "unsupervised learning") Sort by: PublicationDate	34840

2) Embase Query

No.	Query	Items found
#25	#23 AND #24	688
#24	[embase]/lim NOT [medline]/lim	9568801
#23	#20 AND #21 AND ([english]/lim OR [french]/lim OR [italian]/lim OR [spanish]/lim)	1423
#22	#20 AND #21	1433
#21	omic*:ti,ab OR 'machine learning':ti,ab OR 'personalized medicine':ti,ab OR 'personalised medicine':ti,ab	59092
#20	#4 AND #10 AND #16 AND #19	4830
#19	#17 OR #18	6287177
#18	validation:ti,ab OR validity:ti,ab OR validated:ti,ab OR 'cross validation':ti,ab OR 'cross validated':ti,ab OR test*:ti,ab OR 'clinical utility*':ti,ab OR accuracy:ti,ab OR robustness:ti,ab OR reliability*:ti,ab OR sensitivity:ti,ab OR specificity:ti,ab OR benchmark*:ti,ab OR bias:ti,ab OR 'cross study:ti,ab' OR 'cross studies':ti,ab	6150811
#17	'validation study'/exp OR 'reliability'/exp OR 'sensitivity and specificity'/exp OR 'benchmarking'/exp	580344
#16	#14 OR #15	1174400
#15	omic*:ti,ab OR 'omic based':ti,ab OR 'multi omic*':ti,ab OR genomic*:ti,ab OR transcriptomic*:ti,ab OR proteomic*:ti,ab OR metabolomic*:ti,ab OR lipidomic*:ti,ab OR epigenomic*:ti,ab OR microarray:ti,ab OR 'rna seq':ti,ab OR 'mass spectrometr*':ti,ab	852339
#14	#11 OR #12 OR #13	758269
#13	'mass spectrometry'/exp	455591
#12	'microarray analysis'/exp	68369
#11	'omics'/exp OR 'genomics'/exp OR 'epigenetics'/exp	299009
#10	#5 OR #6 OR #7 OR #8 OR #9	5000844
#9	'individualized medicine':ti,ab OR 'individualised medicine':ti,ab OR 'individualized therapy':ti,ab OR 'individualised therapy':ti,ab	3459
#8	'personalised medicine':ti,ab	1713
#7	'personalized medicine':ti,ab	13669
#6	'stratified medicine':ti,ab OR cluster*:ti,ab OR 'sub group*':ti,ab OR subgroup*:ti,ab OR biomarker*:ti,ab OR diagnos*:ti,ab OR prognos*:ti,ab OR 'precision medicine':ti,ab	4904827
#5	'biological marker'/exp OR 'personalized medicine'/exp	330768
#4	#1 OR #2 OR #3	200079
#3	'machine learning'/exp	193633
#2	'statistical learning'/exp	46
#1	'machine learning':ti,ab OR 'statistical learning':ti,ab OR 'supervised learning':ti,ab OR 'unsupervised learning':ti,ab	34557

3) Web of Science Query

No.	Query	Items found
# 7	#6 AND #5 Indexes=SCI-EXPANDED, SSCI, A&HCI, CPCI-S, CPCI-SSH, ESCI Timespan=All years	193
# 6	TITLE: ((omic* OR "machine learning" OR "personalized medicine" OR "personalised Medicine")) Indexes=SCI-EXPANDED, SSCI, A&HCI, CPCI-S, CPCI-SSH, ESCI Timespan=All years	33,111
# 5	#4 AND #3 AND #2 AND #1 Indexes=SCI-EXPANDED, SSCI, A&HCI, CPCI-S, CPCI-SSH, ESCI Timespan=All years	1,075
# 4	TOPIC: ((validation OR validity OR validated OR "cross validation" OR "cross validated" OR "clinical utility*" OR accuracy OR robustness OR reliability* OR sensitivity OR specificity OR benchmark* OR bias OR "cross study" OR "cross studies")) Indexes=SCI-EXPANDED, SSCI, A&HCI, CPCI-S, CPCI-SSH, ESCI Timespan=All years	5,015,557
# 3	TOPIC: ((Omic* OR "omic based" OR "multi omic*" OR genomic* OR transcriptomic* OR proteomic* OR metabolomic* OR lipidomic* OR epigenomic* OR microarray OR "RNA seq" OR "mass spectrometr*")) Indexes=SCI-EXPANDED, SSCI, A&HCI, CPCI-S, CPCI-SSH, ESCI Timespan=All years	1,024,118
# 2	TOPIC: ((" stratified medicine" OR cluster* OR "sub group*" OR Subgroup* OR biomarker* OR diagnos* OR prognos* OR "precision medicine" OR "personalized medicine" OR "personalised medicine" OR "individualized Medicine" OR "individualised Medicine" OR "individualized therapy" OR "individualised therapy")) Indexes=SCI-EXPANDED, SSCI, A&HCI, CPCI-S, CPCI-SSH, ESCI Timespan=All years	4,137,042
# 1	TOPIC: (("Machine learning" OR "statistical learning" OR "supervised learning" OR "unsupervised learning")) Indexes=SCI-EXPANDED, SSCI, A&HCI, CPCI-S, CPCI-SSH, ESCI Timespan=All years	134,956

Online Supplementary file 3 – Data extraction form

Data items extracted from each processed article during the full-text scoping review, and associated qualifications for each item.

Item	Qualifications
Authors	
Title	
Journal	
Volume	
Issue	(if applicable)
Pages	(if applicable)
Year	
Location	
URL / DOI	
Type of publication	<ul style="list-style-type: none"> • Research article • Meeting abstract • Review
Study population and sample size	(if applicable)
Methodology / Study Design	<ul style="list-style-type: none"> • Case-control study • Cases only stratification study <p>(+ further qualification, e.g. treatment response prediction, tumor subtype categorization, recurrence/relapse prediction, survival prediction, tissue-of-origin prediction)</p>
Outcome assessment	<ul style="list-style-type: none"> • Performance measures (e.g. accuracy, sensitivity, specificity, Kohen's Kappa, F-score, AUC) • Validation scheme (cross-validation approach, external validation approach, single cohort or multiple cohorts)
Generic machine learning category	<ul style="list-style-type: none"> • Supervised learning • Unsupervised learning • Other / mixed approaches
Name of specific machine learning approach	(if applicable)

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Main results / key findings that relate to the research question	short description
---	-------------------

For peer review only

Online Supplementary file 1 – PRISMA-ScR Checklist

Preferred Reporting Items for Systematic reviews and Meta-Analyses extension for Scoping Reviews (PRISMA-ScR) Checklist (17).

SECTION	ITEM	PRISMA-ScR CHECKLIST ITEM	REPORTED ON PAGE #
TITLE			
Title	1	Identify the report as a scoping review.	2
ABSTRACT			
Structured summary	2	Provide a structured summary that includes (as applicable): background, objectives, eligibility criteria, sources of evidence, charting methods, results, and conclusions that relate to the review questions and objectives.	2
INTRODUCTION			
Rationale	3	Describe the rationale for the review in the context of what is already known. Explain why the review questions/objectives lend themselves to a scoping review approach.	4
Objectives	4	Provide an explicit statement of the questions and objectives being addressed with reference to their key elements (e.g., population or participants, concepts, and context) or other relevant key elements used to conceptualize the review questions and/or objectives.	4
METHODS			
Protocol and registration	5	Indicate whether a review protocol exists; state if and where it can be accessed (e.g., a Web address); and if available, provide registration information, including the registration number.	(16)
Eligibility criteria	6	Specify characteristics of the sources of evidence used as eligibility criteria (e.g., years considered, language, and publication status), and provide a rationale.	5-6
Information sources*	7	Describe all information sources in the search (e.g., databases with dates of coverage and contact with authors to identify additional sources), as well as the date the most recent search was executed.	5
Search	8	Present the full electronic search strategy for at least 1 database, including any limits used, such that it could be repeated.	5
Selection of sources of evidence†	9	State the process for selecting sources of evidence (i.e., screening and eligibility) included in the scoping review.	6

SECTION	ITEM	PRISMA-ScR CHECKLIST ITEM	REPORTED ON PAGE #
Data charting process‡	10	Describe the methods of charting data from the included sources of evidence (e.g., calibrated forms or forms that have been tested by the team before their use, and whether data charting was done independently or in duplicate) and any processes for obtaining and confirming data from investigators.	6
Data items	11	List and define all variables for which data were sought and any assumptions and simplifications made.	6 (Online Suppl. File 2)
Critical appraisal of individual sources of evidence§	12	If done, provide a rationale for conducting a critical appraisal of included sources of evidence; describe the methods used and how this information was used in any data synthesis (if appropriate).	Click here to enter text.
Synthesis of results	13	Describe the methods of handling and summarizing the data that were charted.	6
RESULTS			
Selection of sources of evidence	14	Give numbers of sources of evidence screened, assessed for eligibility, and included in the review, with reasons for exclusions at each stage, ideally using a flow diagram.	7 (Fig. 2)
Characteristics of sources of evidence	15	For each source of evidence, present characteristics for which data were charted and provide the citations.	7 (Table 1)
Critical appraisal within sources of evidence	16	If done, present data on critical appraisal of included sources of evidence (see item 12).	Click here to enter text.
Results of individual sources of evidence	17	For each included source of evidence, present the relevant data that were charted that relate to the review questions and objectives.	7-11
Synthesis of results	18	Summarize and/or present the charting results as they relate to the review questions and objectives.	7-11
DISCUSSION			
Summary of evidence	19	Summarize the main results (including an overview of concepts, themes, and types of evidence available), link to the review questions and objectives, and consider the relevance to key groups.	11-12
Limitations	20	Discuss the limitations of the scoping review process.	11-12
Conclusions	21	Provide a general interpretation of the results with respect to the review questions and	12-13

SECTION	ITEM	PRISMA-ScR CHECKLIST ITEM	REPORTED ON PAGE #
		objectives, as well as potential implications and/or next steps.	
FUNDING			
Funding	22	Describe sources of funding for the included sources of evidence, as well as sources of funding for the scoping review. Describe the role of the funders of the scoping review.	20

For peer review only

BMJ Open

Biomarker discovery studies for patient stratification using machine learning analysis of omics data: a scoping review

Journal:	<i>BMJ Open</i>
Manuscript ID	bmjopen-2021-053674.R1
Article Type:	Original research
Date Submitted by the Author:	12-Aug-2021
Complete List of Authors:	Glaab, Enrico; University of Luxembourg, Luxembourg Centre for Systems Biomedicine Rauschenberger, Armin; University of Luxembourg, Luxembourg Centre for Systems Biomedicine Banzi, Rita; Mario Negri Institute for Pharmacological Research, Center for Health Regulatory Policies Gerardi, Chiara; Mario Negri Institute for Pharmacological Research, Center for Health Regulatory Policies Garcia, Paula; ECRIN, European Clinical Research Infrastructure Network Demotes, Jacques; ECRIN, European Clinical Research Infrastructure Network
Primary Subject Heading:	Patient-centred medicine
Secondary Subject Heading:	Diagnostics, Research methods
Keywords:	BIOTECHNOLOGY & BIOINFORMATICS, NATURAL SCIENCE DISCIPLINES, STATISTICS & RESEARCH METHODS

SCHOLARONE™
Manuscripts



I, the Submitting Author has the right to grant and does grant on behalf of all authors of the Work (as defined in the below author licence), an exclusive licence and/or a non-exclusive licence for contributions from authors who are: i) UK Crown employees; ii) where BMJ has agreed a CC-BY licence shall apply, and/or iii) in accordance with the terms applicable for US Federal Government officers or employees acting as part of their official duties; on a worldwide, perpetual, irrevocable, royalty-free basis to BMJ Publishing Group Ltd ("BMJ") its licensees and where the relevant Journal is co-owned by BMJ to the co-owners of the Journal, to publish the Work in this journal and any other BMJ products and to exploit all rights, as set out in our [licence](#).

The Submitting Author accepts and understands that any supply made under these terms is made by BMJ to the Submitting Author unless you are acting as an employee on behalf of your employer or a postgraduate student of an affiliated institution which is paying any applicable article publishing charge ("APC") for Open Access articles. Where the Submitting Author wishes to make the Work available on an Open Access basis (and intends to pay the relevant APC), the terms of reuse of such Open Access shall be governed by a Creative Commons licence – details of these licences and which [Creative Commons](#) licence will apply to this Work are set out in our licence referred to above.

Other than as permitted in any relevant BMJ Author's Self Archiving Policies, I confirm this Work has not been accepted for publication elsewhere, is not being considered for publication elsewhere and does not duplicate material already published. I confirm all authors consent to publication of this Work and authorise the granting of this licence.

Title

Biomarker discovery studies for patient stratification using machine learning analysis of omics data: a scoping review

Authors

Enrico Glaab^{1,*}, Armin Rauschenberger¹, Rita Banzi², Chiara Gerardi², Paula Garcia³, Jacques Demotes-Mainard³, and the PERMIT Group

Affiliations

¹Luxembourg Centre for Systems Biomedicine, University of Luxembourg, Campus Belval, Esch-sur-Alzette, Luxembourg.

²Center for Health Regulatory Policies, Istituto di Ricerche Farmacologiche Mario Negri, Milan, Italy.

³European Clinical Research Infrastructure Network (ECRIN), Paris, France

*Correspondence: enrico.glaab@uni.lu; 14, avenue du Rock'n'Roll, L-4361 Esch-sur-Alzette, Luxembourg; Phone: +352-466644 6186

Word count

Abstract: 300

Main text: 6641

Keywords

Biomarkers, Scoping Review, Omics, Machine Learning, Stratification

Abstract

Objective: To review biomarker discovery studies using omics data for patient stratification which led to clinically validated FDA-cleared tests or laboratory developed tests, in order to identify common characteristics and derive recommendations for future biomarker projects.

Design: Scoping review.

Methods: We searched PubMed, EMBASE and Web of Science to obtain a comprehensive list of articles from the biomedical literature published between January 2000 to July 2021, describing clinically validated biomarker signatures for patient stratification, derived using statistical learning approaches. All documents were screened to retain only peer-reviewed research articles, review articles, or opinion articles, covering supervised and unsupervised machine learning applications for omics-based patient stratification. Two reviewers independently confirmed the eligibility. Disagreements were solved by consensus. We focused the final analysis on omics-based biomarkers which achieved the highest level of validation, i.e., clinical approval of the developed molecular signature as a laboratory developed test or FDA approved tests.

Results: Overall, 352 articles fulfilled the eligibility criteria. The analysis of validated biomarker signatures identified multiple common methodological and practical features that may explain the successful test development and guide future biomarker projects. These include study design choices to ensure sufficient statistical power for model building and external testing, suitable combinations of non-targeted and targeted measurement technologies, the integration of prior biological knowledge, strict filtering and inclusion/exclusion criteria, and the adequacy of statistical and machine learning methods for discovery and validation.

Conclusions: While most clinically validated biomarker models derived from omics data have been developed for personalised oncology, first applications for non-cancer diseases show the potential of multivariate omics biomarker design for other complex disorders. Distinctive characteristics of prior success stories, such as early filtering and robust discovery approaches, continuous improvements in assay design and experimental measurement technology, and rigorous multi-cohort validation approaches, enable the derivation of specific recommendations for future studies.

Strengths and limitations of this study

- This scoping review provides an overview of biomarker discovery studies using machine learning analysis of omics data which have led to clinically validated diagnostic and prognostic tools.
- The review discusses shared characteristics of successful biomarker studies as a guidance for study design, discovery and validation method choices for future projects.
- Data extraction and analysis methods focus on deriving recommendations to optimize the design of prospective studies and improve analysis workflows for retrospective studies.
- The review applied minimum eligibility criteria for sample size and statistical validation, but did not assess the quality of the included studies.

Introduction

Personalised medicine is a rapidly developing area in health care research and practice, which aims at providing more effective and safer therapies tailored to the individual patient, by exploiting subject-specific molecular, clinical and environmental data sources (Box 1).

A central tool in personalised medicine and the focus of this study is the machine learning (ML) analysis of omics profiling data to derive molecular biomarker signatures for disease- or drug-based patient stratification (1). The major goals for ML-based omics biomarker development are to develop more reliable and robust tests for drug response prediction, early diagnosis, differential diagnosis or prognosis of the future clinical disease course (2). Omics-derived biomarker signatures may help to guide treatment decisions, and to focus therapies on the right populations to prevent overtreatment, increase success rates, and reduce costs (3). As a research and information tool, they may enable a better monitoring of disease progression and treatment success, and guide new drug development and discovery (4). In contrast to classical single-molecule biomarker approaches, omics signatures have the potential to provide more sensitive, specific and robust predictions of disease-associated outcomes (5).

However, while biomarker discovery projects using omics data have already led to the successful development of clinically validated diagnostic and prognostic tests (6–15), many biomarker studies are discontinued after early development stages or fail in later clinical validation stages. Dedicated statistical and ML methodologies for omics biomarker discovery and validation have been published, as well as recommendations for study design, implementation and reporting (16,17). The distinctive features and approaches which characterize prior successes in translating omics research findings into clinically validated tests have however not yet been investigated in detail. In order to guide future projects on suitable method choices, there is a need for dedicated studies on the key determinants of previous translational successes in ML-based omics biomarker development.

As part of an EU project on “Personalised Medicine Trials” (PERMIT (18)), funded within the H2020 framework, we have therefore investigated the current methodological practices for personalised medicine, covering ML approaches for omics-based patient stratification as a major focus area. While a broader series of questions was established and examined for the overall scoping review (19), for this manuscript, we focused our analysis on biomarker discovery studies that have led to successful, clinically validated FDA-cleared tests or laboratory developed tests (LDTs), to determine their shared and distinctive characteristics compared to studies with no clinical translation. In particular, we aimed to address the following specific research questions:

- Which omics-derived biomarker discovery studies have led to clinically validated tests for patient stratification (LDTs or FDA-cleared tests)?
- What are the key characteristics shared by successful omics biomarker studies and distinguishing them from previously published biomarker studies which have not yet led to clinically validated tests?
- Which types of model building and validation methods have been used to develop clinically validated biomarker signatures, and what are the lessons learned and recommended workflows?
- Which recommendations and guidelines have been proposed to address common challenges in biomarker development using omics data?

1
2
3 These questions lend themselves to a scoping review, because omics-derived biomarker
4 development is still an evolving field, and a preliminary assessment of the potential scope and size
5 of the available biomedical literature on these topics is required as a first step for further follow-up
6 research. Therefore, the objective of this study was to address the above questions by retrieving and
7 examining the current literature on biomarker discovery and validation studies using omics data and
8 ML approaches.
9

10 11 12 **Methods**

13
14 We conducted a scoping review following the methodological framework suggested by the Joanna
15 Briggs Institute (20). This framework consists of six stages: 1) identifying the research questions, 2)
16 identifying relevant studies, 3) study selection, 4) charting the data, 5) collating, summarising and
17 reporting results, and 6) consultation.
18

19 The scoping review approach was considered most suitable to respond to the broad scope and the
20 evolving nature of the field. Compared to systematic reviews that aim to answer specific questions,
21 scoping reviews present a general overview of the evidence pertaining to a topic and are useful to
22 examine emerging trends, to clarify key concepts and identify gaps (21,22). Before conducting the
23 review, a study protocol was published on the online platform Zenodo (19). Due to the iterative nature
24 of scoping reviews, deviations from the protocol are expected and duly reported when occurred. We
25 used the PRISMA-ScR (Preferred Reporting Items for Systematic reviews and Meta-Analyses
26 extension for Scoping Reviews) checklist to report our results (23) (Online supplementary file 1).
27
28
29
30

31 *Study identification*

32 Relevant studies and documents were identified, balancing feasibility with breadth and
33 comprehensiveness of searches. We searched PubMed, EMBASE and Web of Science (last search
34 date: July 27, 2021) for articles describing supervised or unsupervised ML analyses for biomarker
35 discovery or personalised medicine, including both discovery and validation methods. The relevance
36 of the search methodology was ensured by using a strict multi-stage filtering, considering only
37 articles including at least one relevant search term per category from four categories of keywords
38 (“Personalized medicine / Biomarkers”, “Omics”, “Machine Learning” and “Validation”, covering both
39 synonyms for these terms and closely related keywords, see Fig. 1, illustrating the keyword-based
40 search strategy, and Online Supplementary file 2 for the detailed search queries), and subsequently
41 post-filtering the retrieved articles manually to exclude studies not involving omics-based biomarker
42 research or lacking a description of machine learning and validation analyses (see sections on
43 *Eligibility criteria* and *Study selection*). To cover only relevant scientific content, the scope was limited
44 to journal publications and meeting abstracts from international conferences and workshops, and no
45 other grey literature was included. We restricted inclusion to reports published from January 2000 to
46 July 2021 (covering also “online first” articles with official publication date in the future) in English,
47 French, Spanish, Italian and German language. Since to the best of our knowledge, the first clinically
48 validated FDA-cleared omics-derived biomarker signature was published in 2002 (24), only few
49 preliminary discovery studies were expected to have taken place significantly earlier than 2002, and
50 we therefore did not extent the search further backwards in time than January 2000.
51
52
53
54
55
56

57 *Eligibility criteria*

58 We included peer-reviewed methodology articles, review articles, opinion articles on supervised and
59 unsupervised ML methods for omics disease prediction and stratification and associated statistical
60 cross-validation and multi-cohort validation methods (addressing accuracy, robustness, and clinical

1
2
3 relevance). Only approaches tested on real-world biomedical omics data were reviewed, while
4 studies relying purely on simulated data or were excluded. We also excluded papers on biomarker
5 methods without a demonstrated biomedical application, and those with insufficient sample size (i.e.,
6 removing studies covering less than 50 samples per group for the main conditions studied, unless a
7 dedicated power calculation was presented) or statistical validation (i.e., lack of clear descriptions of
8 cross-validation or external testing methodology, performance metrics and test statistics). These
9 exclusion criteria were not specified in the generic review protocol, but they were agreed among the
10 authors prior to the screening process.
11

12
13 To cover both data from original research papers and prior systematic reviews, we extracted
14 information from three main article types: (1) applied research papers, (2) methodology articles with
15 demonstrated applications, and (3) review articles on methods, applications and validation
16 approaches.
17

18 Apart from these inclusion and exclusion criteria, for the final result presentation, the statistical
19 investigations covered all selected articles, whereas the detailed discussion of study characteristics
20 focused on the studies that led to clinically validated biomarker signatures tested on multiple cohorts
21 with large sample sizes (i.e., studies using a power calculation to demonstrate the adequacy of the
22 chosen sample sizes, or covering hundreds or thousands samples per studied subject group).
23
24
25

26 *Study selection*

27
28 We exported the references retrieved from the searches into the online tool Rayyan (25). Duplicates
29 were removed automatically using the reference manager Endnote X9 (Clarivate Analytics,
30 Philadelphia, United States) and manually by the reviewers. One reviewer loaded the retrieved
31 records into the online screening tool Rayyan (25), and two reviewers confirmed the eligibility
32 independently by covering both the screening for all records and the full-text review for the articles
33 pre-selected by the screening. Disagreements were solved by consensus.
34
35
36

37 *Charting the data and synthesis of results*

38
39 We designed a data extraction form using Excel (Online supplementary file 3). General study
40 characteristics extracted covered author names, title, citation, type of publication (e.g., journal article,
41 meeting abstract), study population and sample size (if applicable), methodology/study design, and
42 outcome measures (if applicable). Specific items associated with the topic of the scoping review
43 included the study type (e.g., Case-control study, differential diagnosis study, prognostic study,
44 review – methods, review – applications, review – validation); the article type (journal or conference
45 article), the generic ML domain (e.g., supervised/unsupervised); and the name of specific
46 approaches for outcome prediction and for validation. Moreover, to capture key findings related to
47 the review questions, relevant sentences were extracted from each reviewed article, and if needed,
48 complemented by a brief explanatory remark, and by writing out abbreviations used in the original
49 text.
50
51

52 The reviewers piloted the data extraction form using five records from the retrieved article collection.
53 Two reviewers (EG, AR) working independently extracted the data from the included articles. In the
54 case of disagreements, consensus was obtained by discussion.
55

56 In the final full-text review stage, the pre-selected articles were grouped by topic, categorizing articles
57 into applied vs. methodological studies, supervised vs. unsupervised analyses, and assigning
58 algorithm type identifiers to each article (review articles and papers on validation methodologies
59 were considered as separate categories without a specific algorithm type assignment). The full-text
60

1
2
3 review and categorization of articles into different publication types was done through independent
4 manual inspection by the two reviewers.
5

6 While the information on sample sizes and validation methods was documented as part of the data
7 extraction, it was not within the remit of this scoping review to assess the methodological quality of
8 individual studies included in the analysis.
9

10 11 12 *Consultation exercise*

13 The members of the PERMIT consortium, associated partners, and the PERMIT project Scientific
14 Advisory Board discussed the preliminary findings of the scoping review in a 2-hour online workshop.
15
16

17 18 *Patient and public involvement*

19 The European Patients' Forum is a member of PERMIT project. Although not directly involved in the
20 conduction of the scoping review, they received the draft review protocol for collecting comments
21 and feedback.
22
23
24
25
26
27

28 **Results**

29 *Study selection and general characteristics of reports*

30 We retrieved 1563 abstracts from the literature search. After the removal of duplicates, we screened
31 the remaining 1475 abstracts for eligibility. 619 records were excluded, while 856 abstracts were
32 retained for the full-text assessment. Finally, we included 352 articles that passed all filtering criteria
33 in the data extraction and analysis (see flow chart in Fig. 2).
34
35
36

37 The full-text article review revealed that many studies did not meet the pre-defined inclusion criteria:
38 371 articles (43%) were removed because of an insufficient sample size, and 105 further articles
39 (12%) were excluded because they provided insufficient details on the validation results or
40 methodology (see Fig. 2). This shows that the challenges of recruiting an adequate number of
41 participants per study group or conducting sufficient omics profiling experiments for robust model
42 building and validation are not met in a large proportion of omics biomarker studies. Moreover, many
43 studies lack adequate documentation for the study design and validation.
44

45 For the selected articles that cover primary research on omics biomarker studies, the majority (77%)
46 rely entirely on an internal validation involving data from only a single cohort, whereas studies that
47 use an external validation on an independent cohort are still underrepresented (only 12% of articles
48 describe both an internal cross-validation and an external cohort validation, and an additional 11%
49 include an external validation, but do not report internal cross-validation results). However, when
50 comparing the numbers of published studies over different periods of time during the past 20 years,
51 the relative proportion of studies including an external validation has increased in recent years (see
52 Fig. 3), suggesting a growing recognition of the importance of independent, multi-cohort validation.
53

54
55 Next, we investigated the countries of origin for the selected articles, showing that the United States
56 of America (USA) are contributing the largest proportion of validated biomarker studies (28%),
57 followed by China (18%), Canada (5%), Germany (4%), and the United Kingdom and India (both
58 3%; see also Fig. 4, providing a map visualization of the country statistics). These country
59 representations show limited correlation with population sizes and may largely reflect worldwide
60 variation in relative biomedical research productivity reviewed in previous study (26). Since the most

1
2
3 prolific countries in the development of molecular diagnostics have already set up policies and
4 regulations for omics- and ML-based *in vitro* diagnostics and medical devices (e.g., see the life cycle
5 regulation of AI- and ML-based software devices in the USA (27)), they may also provide a role
6 model for countries still in the process of establishing similar regulatory frameworks.
7

8
9 When inspecting the representation of study design types in the filtered article collection, the great
10 majority of documents described diagnostic studies (67%), prognostic and survival prediction studies
11 were covered in 8% of articles, and studies examining therapy or drug response in 7% (see Fig. 5).
12 Apart from this, 13% of articles were reviews on methodologies and applications in the field, and 5%
13 of articles described other rare study types (e.g. tissue-of-origin prediction studies or combinations
14 of different study types).
15

16 Since a detailed discussion of all filtered articles is not within the scope of the present review, in the
17 following, we focus on reviewing representative omics biomarker studies which achieved the highest
18 validation level, i.e., clinical approval of the developed molecular signature as an LDT or FDA
19 approved test (see the overview of studies in Table 1 and the FDA web-site (28)). We investigate
20 the shared features of these successful studies, examine how they address common shortcomings
21 and missing features of other reviewed studies, and summarize the lessons learned.
22
23
24
25
26

27 *Success stories in omics-based biomarker signature development*

30 *Cancer approved omics-derived diagnostic tests (9 studies)*

31
32 The first and most well-known omics-derived molecular test to receive FDA clearance was
33 *MammaPrint*[®], a prognostic signature using the RNA expression activity of 70 genes to estimate the
34 risk for distant tumours metastasis and recurrence in early-stage breast cancer patients (6,24,29–
35 32). This test was developed at the Netherlands Cancer Institute, using DNA microarray analysis to
36 investigate primary breast tumours of 117 patients. Supervised ML was applied to the resulting data
37 to identify a highly predictive gene signature for a short interval to distant metastases in lymph node
38 negative patients (24).
39
40

41 A distinctive feature of the development approach behind this signature in comparison to other
42 reviewed studies was the multi-stage filtering and cross-validation strategy used in the initial
43 discovery study, which may explain the repeated confirmation of the signature in later validation
44 studies (6,29–32). From 25k genes represented on the DNA microarrays, only those significantly
45 regulated in more than 3 tumours out 78 sporadic lymph-node negative patients were preselected,
46 and further filtered by retaining only the genes with a minimum absolute correlation with the disease
47 outcome of 0.3. The resulting list of 231 genes, rank-ordered by absolute correlation, was
48 investigated by sequentially adding the next top 5 genes from the list to a candidate ML classifier
49 and evaluating its performance by leave-one-out cross-validation (LOOCV). This procedure was
50 repeated as long as the estimated accuracy of the classifier improved, providing a final candidate
51 signature of 70 genes. The final signature was validated on multiple independent test sets, including
52 a set of 19 external samples in the original study and several additional validations on independent
53 cohorts in follow-up studies (6,29–32).
54
55

56 The *MammaPrint* signature provided the role model for the subsequent development of a similar
57 prognostic test for colon cancer, *ColoPrint*[®] (33–38). This test aims at detecting the approx. 20% of
58 patients with stage II colon cancer expected to experience a relapse and develop distant metastases.
59 It uses an 18-gene expression signature, developed by analysing DNA microarray data in a similar
60 manner to the *MammaPrint* approach. The diagnostic approach has been commercialized as an

1
2
3 LDT to assist physicians in selecting treatment options for colon cancer patients. Similar to
4 *MammaPrint*, the signature development was characterized by extensive discovery and validation
5 studies, which involved multiple statistical reproducibility, stability and precision analyses for
6 independent, large-scale patient cohorts (39).
7

8 Another widely used cancer-related LDT, which received FDA clearance in 2013, is the *Prosigna*[®]
9 *Breast Cancer Prognostic Gene Signature Assay*, previously called *PAM50* test (40–44). This assay
10 assesses mRNA expression for a signature of 58 genes (50 target genes + 8 endogenous control
11 genes) to predict the risk of distant recurrence for hormone-receptor-positive breast cancer between
12 5 to 10 years after diagnosis (prerequisites are that the patients have been treated with hormonal
13 therapy and surgery, and are stage I or stage II lymph-node negative, or in stage II with one to three
14 positive nodes). The test development started with a microarray discovery study and involved a
15 multistage filtering, using consecutive applications of statistical tests and cross-validation to propose
16 a subset of candidate gene markers (45). The authors compared the reproducibility of classification
17 scores obtained with these markers for three centroid-based prediction methods to ensure the
18 robustness of the methodology. By further developing the approach into a more sensitive PCR-based
19 test, and later into an assay using the NanoString nCounter Dx Analysis System, the predictive
20 performance was improved in a step-wise fashion. The original discovery study was characterized
21 by significantly larger sample sizes than the majority of reviewed biomarker studies, with a training
22 set of 189 samples, test sets of 761 patients evaluated for prognosis, and 133 patients evaluated for
23 prediction of pathologic complete response to treatment with taxane and anthracycline. These study
24 design features in combination with multi-stage filtering and validation approaches, and improved
25 measurement technology during the course of the study, may explain the successful progression of
26 the *PAM50* test to FDA clearance.
27
28
29

30 Among the LDTs for breast cancer prognosis, *Oncotype DX*[®] is a further test commonly used in
31 clinical practice (8,46–49). The underlying gene signature consists of 16 cancer-associated genes
32 and 5 reference genes, and is therefore often also referred to as ‘21-gene assay’. Its main application
33 is to predict risk of recurrence in oestrogen-receptor positive tumours. The relevance of this
34 prognostic tool for treatment selection may be explained by the strong association of the provided
35 recurrence score with the probability of positive treatment response to chemotherapy (50). *Oncotype*
36 *DX* was developed using a consecutive refinement procedure, starting with the RT-PCR assessment
37 of 250 candidate genes across 447 patients from three distinct studies to identify the 21-gene
38 signature after multiple filtering steps. A recurrence score algorithm built using the signature as input
39 was clinically validated on 668 independent patients (51). The selection of the 16 cancer-related
40 genes included in the assay involved scoring the performance of the candidate features in all three
41 studies and the consistency of the primer/probe performance in the assay (52). Thus, particular
42 strengths of the development process for this LDT include the consideration of both technical
43 robustness and statistical robustness of the assay across distinct cohorts. However, an independent
44 comparative clinical validation of *Oncotype DX* and the *PAM50* signature for estimating the likelihood
45 of distant recurrence in ER-positive, node-negative, post-menopausal breast cancer patients treated
46 with endocrine therapy suggested that the *PAM50* signature provided more prognostic information
47 than *Oncotype DX* (53).
48
49
50

51 While the first validated omics biomarker signatures were developed for breast cancer, similar
52 diagnostic and prognostic tools have followed for other cancer types. One of these is the *Decipher*[®]
53 *Prostate Cancer Test* (9,54–58), which differs from other omics-derived diagnostic tools by being
54 provided together with a software platform and database, the *Decipher Genomic Resource*
55 *Information Database (GRID)*, that captures 1.4 million expression markers per patient to facilitate
56 personalised care. The test itself uses 22 preselected RNAs to predict clinical metastasis and
57 cancer-specific mortality for patients who have undergone radical prostatectomy. An initial discovery
58 study by the Mayo Clinic (Rochester, MN, USA) investigated a cohort of 545 such patients, split into
59 a training (n = 359) and a validation cohort (n = 186). Similar to other LDTs, the discovery started
60

1
2
3 with a genome-wide profiling and used both statistical and ML analyses for filtering. First, *t*-tests
4 were applied (reduction from 1.4 mil. to 18,902 differentially expressed RNAs), then regularized
5 logistic regression (reduction to 43 candidate markers), and finally a random forest-based feature
6 selection (reduction to final set of 22 RNAs). Apart from testing the signature in the validation cohort,
7 further external validations were performed in subsequent studies (9,54–58). Overall, distinctive
8 strengths of the used approach include the improved interpretability of the test results through
9 supporting analyses on the GRID platform, and the robustness of the discovery and validation
10 approach, involving large sample sizes and several complementary statistical and ML assessments.
11

12
13 While most diagnostic tests in oncology have been designed for specific cancer types, a dedicated
14 LDT has also been developed for cancers of unknown or uncertain diagnosis. The Cancer Type ID®
15 test by bioTheranostics distinguishes between 50 different tumour types using a 92-gene RT-PCR
16 expression measurement signature (15,59–61). This signature was derived from analyses of a
17 microarray data collection covering 446 frozen tumour samples and 112 formalin-fixed, paraffin-
18 embedded (FFPE) samples of both primary and metastatic tumours. Modelling steps involved *k*-
19 nearest neighbour clustering and classification, and a genetic algorithm to explore the search space
20 of possible feature subset selections. After successful cross-validation (84% accuracy) and external
21 validation (82% accuracy on 112 independent FFPE samples), the microarray-based signature was
22 further developed to use more sensitive RT-PCR measurements. Testing the new approach on an
23 independent validation set provided an increased accuracy (87%). Distinctive characteristics of the
24 development process that may have contributed to the positive validation include the efficient and
25 extensive exploration of the search space of possible gene subset selections via a genetic algorithm,
26 the large sample sizes used for discovery and validation, and the transfer of the assay from
27 microarrays to the more sensitive RT-PCR platform.
28

29
30 The first omics-derived biomarker signatures addressed only the most frequent cancer types, but
31 more recent applications in oncology focus on the diagnosis of less common malignancies, such as
32 thyroid cancer. Typically, deciding whether a thyroid nodule is benign or cancerous is possible via a
33 fine needle aspiration (FNA) biopsy, without requiring more complex measurements or analyses.
34 However, while direct FNA-based diagnosis is feasible in most cases, indeterminate results can
35 occur (62). To help prevent unnecessary surgeries for the corresponding patients, a molecular
36 signature and LDT known as the Afirma™ Gene Expression Classifier (GEC) has been developed
37 to discriminate benign from cancerous thyroid nodules (62–67). The original discovery study behind
38 the GEC signature used mRNA expression analysis in 315 thyroid nodules, covering 178
39 retrospective surgical tissues and 137 prospectively collected FNA samples. Two ML classifiers were
40 trained separately on surgical tissues and FNAs, assessing the test set performance on 48
41 independent, prospective FNA samples (50% of which had indeterminate cytopathology).
42 Discriminative features were selected using a linear modelling approach implemented in the software
43 Limma, and a linear support vector machine was applied for model building and performance
44 estimation via 30-fold cross-validation (CV). The successful cross-validation results were confirmed
45 on multiple distinct cohorts (62,65–68). While the internal validation used in the initial study cannot
46 address cohort-specific biases, the combined use of established feature selection and modelling
47 approaches, and the subsequent external validation across multiple cohorts with large sample sizes
48 may account for the successful translation of this signature.
49

50
51
52 Most omics-based diagnostic tests identified in our study rely purely on gene expression profiling
53 data. However, more recently, first multi-omics signatures for diagnostic purposes have been
54 developed. One of the first LDTs that integrated information from both RNA and DNA sequencing
55 was the FoundationOne® Heme assay (14,69–71). This assay aims to detect hematologic
56 malignancies, sarcomas, pediatric malignancies, or solid tumours (including among others
57 leukaemias, myelodysplastic syndromes, myeloproliferative neoplasms, lymphomas, multiple
58 myeloma, Ewing sarcoma, Leiomyosarcoma, and paediatric tumours). The test identifies four types
59 of genomic alterations (base substitutions, insertions and deletions, copy number alterations,
60

1
2
3 rearrangements) and reports microsatellite instability and tumour mutational burden to facilitate
4 clinical decision making. This approach was originally developed and evaluated using reference
5 samples of pooled cell lines in order to model the main characteristics that determine the test
6 accuracy, including mutant allele frequency, indel length and amplitude of copy change (69). A first
7 validation using 249 independent FFPE cancer samples, which had already been characterized by
8 established assays, confirmed the accuracy of the test. External validation studies on independent
9 cohorts corroborated the utility of the test for further diagnostic applications (14,72). The study results
10 highlight the potential of integrating diverse biological data sources in order to obtain more robust
11 and reliable predictions, a strategy that may be promising in particular for complex disorders that
12 involve very heterogeneous phenotypes.
13
14

15 A common limitation of genomic profiling approaches for diagnostic testing is that most analyses
16 have to be performed in centralized specialty laboratories, which limits a wider use and results in
17 long waiting times. To address this shortcoming, the Elio™ Tissue Complete assay, an *in vitro*
18 diagnostic test cleared in 2020 by the FDA for assessing somatic mutations and tumour mutation
19 burden (TMB) in solid tumours, has been developed as an integrated DNA-to-report approach to
20 enable a decentralized evaluation in all diagnostic labs with next generation sequencing (NGS)
21 technology (73). The analytical performance of the test was assessed by comparing it with the
22 FoundationOne test (see above) using a concordance analysis on 147 tumour specimens. It
23 provided a positive percent agreement (PPA) above 95% for single nucleotide variants (SNVs) and
24 insertions/deletions, and 80-83% PPA for copy number alterations and gene translocations (73). The
25 test has recently also been applied to investigate the response to immune checkpoint inhibitors (ICI)
26 in metastatic renal cell carcinoma (mRCC), using a retrospective evaluation of SNVs, TMB,
27 microsatellite status and genomic status of antigen presentation genes (74). While no correlation
28 between treatment response and TMB was observed, one third of patients with progressive disease
29 following ICI therapy displayed loss of heterozygosity of major histocompatibility complex class I
30 genes (LOH-MHC) vs. 6% of disease control patients, suggesting that loss of antigen presentation
31 may restrict ICI response (74). In summary, the Elio Tissue Complete assay provides an example of
32 how integrating NGS analyses with bioinformatics in a combined DNA-to-report approach could help
33 to broaden the access to genomic diagnostics for both clinical and research applications.
34
35
36
37
38

39 *Non-cancer approved omics-derived diagnostic tests (4 Studies)*

40
41 While most clinically approved omics-derived diagnostic tests have been developed in the field of
42 oncology, one of the first LDTs that received FDA clearance for a non-cancer disease was the
43 AlloMap® Heart test (13,75–77). It uses a gene expression signature of 11 target genes and 9 control
44 genes in peripheral blood from heart transplant recipients to estimate the risk for acute cellular
45 cardiac allograft rejection. The development process involved statistical analyses of leukocyte
46 microarray profiling data from 285 samples, and subsequent RT-PCR validation and bioinformatics
47 post-processing (13). Prior knowledge from database and literature mining was included in the
48 analysis by mapping the data to known alloimmune pathways. This allowed the researchers to
49 narrow down 252 candidate marker genes. An RT-PCR validation on 145 samples confirmed 68 of
50 these candidate genes, which distinguished rejection samples from quiescent samples according to
51 a T-test ($p < 0.01$). Six genes were eliminated due to significant variation in gene expression with
52 sample processing time. Next, the investigators averaged correlated gene expression levels to
53 create robust meta-level features, called 'metagenes', and added 20 of these features as new
54 variables. A linear discriminant analysis was applied, providing a prediction model using four
55 individual genes and three metagenes, which aggregate information from 11 original genes. Finally,
56 bootstrap validation procedures and external test set validations were performed to confirm the
57 accuracy of this signature. Overall, distinctive aspects of the development approach for the AlloMap
58 signature include the knowledge-based gene discovery, a comprehensive RT-PCR validation of
59 candidate genes, and the robust bootstrap and external validation analyses.
60

1
2
3 The first clinically validated LDT for a cardiovascular indication derived from omics data was the
4 Corus® CAD test, developed to identify coronary artery disease (CAD) in stable non-diabetic patients
5 (11,78–81). In contrast to most other omics-based tests, Corus CAD is not a pure molecular
6 signature test, but takes the clinical covariates gender and age into account. The initial discovery
7 study used a retrospective microarray analysis of blood samples from 195 diabetic and non-diabetic
8 patients from the Duke University CATHGEN registry. After ranking the studied genes by the
9 statistical significance of group differences and prior biological knowledge on their disease
10 relevance, 88 genes were selected for RT-PCR validation. Because diabetes status as a clinical
11 covariate was significantly associated with the observed gene expression alterations, and the
12 identified CAD-associated genes did not overlap between diabetic and non-diabetic patients, the
13 authors decided to limit follow-up work to non-diabetic patients. In a prospective clinical trial,
14 microarray profiling was conducted on blood samples from 198 patients, and top-ranked genes were
15 further validated using RT-PCR for 640 blood samples. After multiple filtering steps, taking into
16 account statistical significance in T-tests, biological relevance, gene correlation clustering and cell-
17 type analyses, a final signature of 23 genes was derived, composed of 20 CAD-associated genes
18 and 3 reference genes (82). To maximize the predictive performance, the final prediction algorithm
19 was optimized to adjust for differences associated with age and gender. Compared to most other
20 reviewed studies, the Corus CAD approach stands out by taking clinical covariates into account in
21 the final prediction model, including an intermediate critical review and adjustment of the inclusion
22 criteria (limiting the focus to nondiabetic patients), and integrating complementary filtering and
23 validation analyses on large sample sizes.
24
25
26

27 For inflammatory diseases, a first omics-derived signature recently received approval for measuring
28 rheumatoid arthritis (RA) inflammatory disease activity, the Vectra® DA multi-biomarker test (83–87).
29 It uses blood serum samples and multi-spot 96-well immunoassay plates to assess serum
30 concentrations of 12 protein biomarkers associated with the pathobiology of RA. The original Vectra
31 DA score, which combines these measurements into a composite score between 1 and 100, was
32 assessed via multivariate regression and displayed a high predictive power in estimating a standard
33 RA score, the Disease Activity Score in 28 joints using the C-reactive protein level (DAS28-CRP), in
34 both seropositive (AUC 0.77, $P < 0.001$) and seronegative (AUC 0.70, $P < 0.001$) patients (87). This
35 score was later adjusted for age, gender and adiposity (based on leptin concentration), and validated
36 in two cohorts against DAS28-CRP as a prognostic test for radiographic progression during the next
37 year. The results showed that the new adjusted score was the most accurate, independent predictor
38 of progression, with the rate of progression increasing from $< 2\%$ in the low (1-29) adjusted score
39 category to 16% in the high (45-100) category (85). Overall, the Vectra DA approach illustrates the
40 utility of omics-based biomarker signatures for prognostic applications in inflammatory disorders,
41 and further highlights the benefit of integrating omics signatures with information from clinical
42 covariates.
43
44
45

46 For neurodegenerative disorders, clinically approved diagnostic and prognostic omics-derived tests
47 are still lacking. However, recently the Helix® Genetic Health Risk App for Late-onset Alzheimer's
48 Disease (AD) was cleared by the FDA for over-the-counter use. It detects clinically relevant variants
49 in genomic DNA isolated from human saliva of individuals ≥ 18 years in order to report and interpret
50 genetic health risks, and evaluates the information of variants with established genome-wide
51 significant associations to AD. When tested on 99 human saliva samples, the accuracy was 100%
52 with a lower 95% CI bound of 96.3% (88). The approach uses a whole exome sequencing (WES)
53 constituent device, the Helix® Laboratory Platform (89–91), as a qualitative *in vitro* diagnostics
54 approach covering measurements for approximately 20k genes. The Helix Laboratory Platform has
55 received FDA clearance through a new regulatory approval pathway established by the FDA for
56 WES devices (Regulation 21 CFR 866.6000). Due to the generic applicability of the WES profiling
57 assay used by this platform, called Exome+, the assay has also been applied to find statistically
58 significant gene-based associations for several other phenotypes in large-scale cohort studies (89)
59 and to identify carriers of autosomal dominant diseases by population-based genetic screening (91).
60

1
2
3 Thus, the Helix Laboratory Platform provides a first example for a new approval pathway for omics-
4 based diagnostic tests, in which a clinically approved genomic testing device is not anymore linked
5 to a single diagnostic application or a specific disease type. Instead, the market authorization for
6 diagnostic tests is obtained separately from the device and facilitated and accelerated by the prior
7 approval of the constituent measurement device. For the future development of omics-derived
8 biomarker signatures, this may allow researcher to focus on demonstrating the clinical utility of a
9 new signature, while the analytical validity of the underlying testing device has already been
10 established previously.
11
12
13

14 Discussion

15 *Statement of principal findings*

16
17
18 The scoping review of articles on patient stratification using omics data revealed common limitations
19 in the study design for many published biomarker development projects, such as insufficient and
20 imbalanced sample sizes per study group and inadequate validation methods, but also identified
21 multiple studies that have led to validated diagnostic and prognostic tests. These success stories
22 were investigated in more detail to identify common characteristics in the study design, discovery
23 and validation methods, which may have supported the clinical translation of the initial findings. Fig.
24 6 outlines key shared aspects that are possible determinants of the study success and could help to
25 guide future biomarker investigations. In particular, they cover the following main features:
26
27

28 (1) A sample size selection, study group and replicate design that provides adequate statistical
29 power for the ML analyses;

30
31 2) The application of robust statistical filtering and evaluation schemes (including multiple layers of
32 statistical and ML-based feature selection, combined statistical and biological filters, robust
33 validation schemes that involve multiple cross-validation, bootstrapping and external validation
34 analyses, using multiple suitable and complementary performance metrics, and providing
35 information on the statistical variation and confidence intervals for the performance estimates);
36

37 3) Clarity of the study scope and goals (involving clear inclusion and exclusion criteria, primary and
38 secondary outcomes, and decision processes to make necessary adjustments due to new
39 knowledge gained during the project, such as the adjusted inclusion criteria in the Corus CAD study
40 and the progression from non-targeted microarray technology to higher-sensitivity RT-PCR in the
41 case of the Prosigna test and the Cancer Type ID test);
42

43 4) Completeness and reproducibility of the study documentation (covering details on used
44 instruments, parameters and settings, reproducible methods descriptions, and information on data
45 provenance);
46

47 5) Interpretability and biological plausibility of the created predictive models (including explainable
48 and justifiable predictions, human-interpretable model descriptions, and biologically plausible
49 models that agree with the current mechanistic understanding of the studied disorder);
50

51 6) Integration of prior biological knowledge into the predictive feature selection, model building and
52 validation procedures (e.g., using public data on disease-associated molecular pathways and
53 networks; complementary clinical and real-world data, and relevant multi-omics data).
54
55
56
57

58 *Strengths and Limitations*

59 The majority of methodological recommendations derived from the study relate to the early planning
60 and study design for biomarker discovery projects, involving considerations associated with the

1
2
3 choice of the study group, sampling and blocking design, the measurement technology, and the input
4 and output variables (16, 17). These recommendations are therefore mainly applicable to prospective
5 studies. For retrospective biomarker investigations of already collected data, the suggestions derived
6 from the review are limited to guidance on improving analysis workflows, e.g. for filtering and
7 evaluation analyses, the integration of prior knowledge from multi-omics data and public annotation
8 databases, and the choice of robust and interpretable modelling approaches for the generation of
9 biologically plausible and reproducible prediction models. While the focus of the review on studies
10 that have already led to validated biomarker models and that fulfil minimum requirements for sample
11 size and statistical model assessment helps to ensure the quality of the selected articles, no further
12 quality evaluation was performed. Finally, more recent methodological developments in the machine
13 learning and cross-validation analysis of omics data, such as meta-learning (92) and bolstered cross-
14 validation (93), have only limited coverage among the articles that passed the eligibility criteria, and
15 will therefore require further dedicated study in the future.
16
17
18
19

20 *Discussing important differences in results*

21
22 Previous reviews of ML approaches using omics data for patient stratification have focused on
23 domain-specific analyses for specific types of diseases, or specific types of ML methodologies (94–
24 102). By contrast, this scoping review focuses on disease-agnostic workflows with generic
25 applicability across complex human disorders involving multifactorial molecular alterations. The
26 coverage of statistical and ML approaches for stratification does not aim to provide a detailed
27 discussion of specific algorithms, statistical methods or scoring metrics, but rather at identifying key
28 determinants of success for generic analysis and validation workflows in biomedical stratification
29 studies. Therefore, the results describe general workflow characteristics that distinguish omics
30 biomarker studies with clinical translation from other studies, and cover associated disease-agnostic
31 recommendations for future studies, whereas method recommendations specific to particular
32 disease types or ML analysis types are covered elsewhere in domain-specific reviews (94–102).
33
34
35
36

37 *Meaning of the study: implications for clinicians and policymakers*

38
39 The previous clinical translation successes in omics-based biomarker development reviewed in this
40 study, which have mostly been achieved in the field of oncology, highlight the potential for developing
41 similar biomarker signatures for further disease indications. In contrast to conventional statistical
42 biomarker discovery approaches, which focus on identifying single-molecule markers, systems-level
43 analysis of omics data using multivariate ML approaches can identify multifactorial signatures which
44 are robust against noise in individual gene or protein measurements, and more biologically insightful
45 by reflecting disease-associated cellular process alterations in a more comprehensive fashion.
46
47

48 This scoping review has identified common characteristics of omics studies which have led to
49 clinically validated diagnostic and prognostic tests. Thus, the conclusions drawn on recommended
50 practices for sample size selection, biological data filtering and ML, and the implementation of
51 adequate validation schemes may help to guide clinical researchers on study design choices and
52 the selection of analysis methodologies. Additionally, the scoping review results can help to raise
53 awareness of common pitfalls, such as issues associated with batch effects, biases, confounding
54 factors, lack of statistical power, and multiple hypothesis testing, and thus contribute to preventing
55 these failure causes in biomarker development. For policymakers and funding bodies, findings on
56 the distinctive characteristics of studies with successful clinical biomarker translation, e.g.
57 concerning the specific requirements for robust cross-validation and external result validation
58 methods, may provide relevant information for the design of public and private funding schemes for
59 biomedical research. Risks in funded research projects may be addressed upfront through
60 appropriate guidelines and regulations for the study design and validation (e.g. recommendations

1
2
3 on power calculations and specific validation and documentation requirements). Finally, the scoping
4 review results can guide clinicians involved in biomarker discovery on how to make better use of
5 available public knowledge and data sources, e.g. cellular pathway and molecular interaction
6 databases, that may allow them to exploit prior knowledge effectively, and create more robust and
7 interpretable biomarker models.
8
9

10 11 *Unanswered questions & future research*

12
13 Since the recommendations and guidelines identified from the reviewed articles are mostly derived
14 from established biomarker discovery and validation approaches, new methodologies and upcoming
15 trends could only be covered to a limited extent and may lead to changed recommendations in the
16 future. In particular, in the reviewed patient stratification studies, some of more recently introduced
17 ML concepts (e.g. transfer learning, distance metric learning, semi-supervised learning, structured
18 machine learning, meta learning, multi-view learning, and generative models), data processing
19 techniques (e.g. new dimension reduction approaches, outlier removal methods, data augmentation
20 techniques), and model validation methods (e.g. bootstrapping or bolstered cross-validation,
21 uncertainty quantification), are still underrepresented among the eligible studies reviewed, and may
22 provide suitable topics for follow-up research.
23
24

25 Overall, while the currently available literature on validated stratification biomarkers already provides
26 ample information on common pitfalls and established practices, the development of widely accepted
27 standard guidelines on methodologies for omics biomarker discovery will require further knowledge
28 exchange and deliberation among stakeholders in the field. In particular, integration of domain-
29 specific expertise in discussions involving clinicians, experimental and data scientists, and regulatory
30 and legal experts is required as a follow-up effort to derive comprehensive methodological guidelines
31 for future biomarker development.
32
33

34 35 36 37 **Acknowledgments**

38
39 The authors thank Vanna Pistotti for her assistance with search strategy development and
40 conduction.
41
42

43 44 **List of Figures**

45
46 Figure 1: Keyword based search strategy for the scoping review

47
48 Figure 2: Study selection flow diagram

49
50 Figure 3: Validation methods used in omics biomarker studies

51
52 Figure 4: Map representation of country statistics for the selected articles

53
54 Figure 5: Representation of study types among the selected articles

55
56 Figure 6: Characteristics of successful omics-based studies
57
58
59
60

Definitions (In boxes)

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Box 1: What is Personalised Medicine?

According to the European Council Conclusion on personalised medicine for patients, personalised medicine is 'a medical model using characterisation of individuals' phenotypes and genotypes (e.g., molecular profiling, medical imaging, lifestyle data) for tailoring the right therapeutic strategy for the right person at the right time, and/or to determine the predisposition to disease and/or to deliver timely and targeted prevention (103).

In the context of the Permit project, we applied the following common operational definition of personalised medicine research: a set of comprehensive methods (methodology, statistics, validation, technology) to be applied in the different phases of the development of a personalised approach to treatment, diagnosis, prognosis, or risk prediction. Ideally, robust and reproducible methods should cover all the steps between the generation of the hypothesis (e.g., a given stratum of patients could better respond to a treatment), its validation and pre-clinical development, and up to the definition of its value in a clinical setting (19).

For peer review only

References

1. Reel PS, Reel S, Pearson E, Trucco E, Jefferson E. Using machine learning approaches for multi-omics data analysis: A review. *Biotechnol Adv.* 2021;49(107739).
2. Goecks J, Jalili V, Heiser LM, Gray JW. How Machine Learning Will Transform Biomedicine. Vol. 181, *Cell.* 2020. p. 92–101.
3. Jiang Y, Wang M. Personalized medicine in oncology: Tailoring the right drug to the right patient. *Biomarkers in Medicine.* 2010.
4. Hopp WJ, Li J, Wang G. Big Data and the Precision Medicine Revolution. *Prod Oper Manag.* 2018;27(9):1647–64.
5. Glaab E. Using prior knowledge from cellular pathways and molecular networks for diagnostic specimen classification. *Brief Bioinform.* 2016;
6. Cardoso F, van't Veer LJ, Bogaerts J, Slaets L, Viale G, Delaloge S, et al. 70-Gene Signature as an Aid to Treatment Decisions in Early-Stage Breast Cancer. *N Engl J Med.* 2016;375(8):717–29.
7. Bachleitner-Hofmann T, Simon I, Salazar R, Tabernero J, Rosenberg R, van der Akker J, et al. Development and Validation of a Robust Molecular Diagnostic Test (COLOPRINT) for Predicting Outcome in Stage II Colon Cancer Patients. *Ann Oncol.* 2012;
8. Knezevic D, Goddard AD, Natraj N, Cherbavaz DB, Clark-Langone KM, Snable J, et al. Analytical validation of the Oncotype DX prostate cancer assay - a clinical RT-PCR assay optimized for prostate needle biopsies. *BMC Genomics.* 2013;14(1).
9. Torres A, Alshalalfa M, Tomlins SA, Erho N, Gibb EA, Chelliserry J, et al. Comprehensive Determination of Prostate Tumor ETS Gene Status in Clinical Samples Using the CLIA Decipher Assay. *J Mol Diagnostics.* 2017;19(3):475–84.
10. Angell TE, Babiarz J, Barth N, Blevins T, Duh Q, Ghossein RA, et al. Clinical validation of the AFIRMA genomic sequencing braf V600E classifier. *Thyroid [Internet].* 2017;27:A50. Available from: <http://www.embase.com/search/results?subaction=viewrecord&from=export&id=L624116485>
11. Ladapo JA, Budoff MJ, Sharp D, Kuo JZ, Huang L, Maniet B, et al. Utility of a Precision Medicine Test in Elderly Adults with Symptoms Suggestive of Coronary Artery Disease. *J Am Geriatr Soc.* 2018;66(2):309–15.
12. Tabari E, Lovejoy AF, Lin H, Bolen CR, Saelee SL, Lefkowitz JP, et al. Molecular characteristics and disease burden metrics determined by next-generation sequencing on circulating tumor DNA correlate with progression free survival in previously untreated diffuse large B-cell lymphoma. *Blood [Internet].* 2019;134. Available from: <http://dx.doi.org/10.1182/blood-2019-123633>
13. Deng MC. The AlloMap™ genomic biomarker story: 10 years after. *Clin Transplant.* 2017;31(3).
14. He J, Abdel-Wahab O, Nahas MK, Wang K, Rampal RK, Intlekofer AM, et al. Integrated genomic DNA/RNA profiling of hematologic malignancies in the clinical setting. *Blood.* 2016;127(24):3004–14.
15. Ma XJ, Patel R, Wang X, Salunga R, Murage J, Desai R, et al. Molecular classification of human cancers using a 92-gene real-time quantitative polymerase chain reaction assay. *Arch Pathol Lab Med.* 2006;130(4):465–73.
16. McShane LM, Cavenagh MM, Lively TG, Eberhard DA, Bigbee WL, Williams PM, et al. Criteria for the use of omics-based predictors in clinical trials. *Nature.* 2013.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

17. Moons KGM, Altman DG, Reitsma JB, Ioannidis JPA, Macaskill P, Steyerberg EW, et al. Transparent reporting of a multivariable prediction model for individual prognosis or diagnosis (TRIPOD): Explanation and elaboration. *Ann Intern Med*. 2015;162(1):W1–73.
18. Banzi R, Gerardi C, Fratelli M, Garcia P, Torres T, Abad JMH, et al. Web-page for the Personalized Medicine Trials (PERMIT) project [Internet]. 2020 [cited 2021 Aug 2]. Available from: <https://permit-eu.org>
19. Banzi R, Gerardi C, Fratelli M, Garcia P, Torres T, Abad JMH, et al. Methodological approaches for personalised medicine: protocol for a series of scoping reviews [Internet]. 10.5281/zenodo.3770937. Available from: <https://zenodo.org/record/3770937>
20. Peters MDJ, Godfrey CM, Khalil H, McInerney P, Parker D, Soares CB. Guidance for conducting systematic scoping reviews. *Int J Evid Based Healthc*. 2015;13(3):141–6.
21. Pham MT, Rajić A, Greig JD, Sargeant JM, Papadopoulos A, McEwen SA. A scoping review of scoping reviews: Advancing the approach and enhancing the consistency. *Res Synth Methods*. 2014;5(4):371–85.
22. Munn Z, Peters MDJ, Stern C, Tufanaru C, McArthur A, Aromataris E. Systematic review or scoping review? Guidance for authors when choosing between a systematic or scoping review approach. *BMC Med Res Methodol*. 2018;18(1):143.
23. Tricco AC, Lillie E, Zarin W, O'Brien KK, Colquhoun H, Levac D, et al. PRISMA extension for scoping reviews (PRISMA-ScR): Checklist and explanation. *Ann Intern Med*. 2018 Oct;169(7):467–73.
24. Van't Veer LJ, Dai H, Van de Vijver MJ, He YD, Hart AAM, Mao M, et al. Gene expression profiling predicts clinical outcome of breast cancer. *Nature*. 2002;415(6871):530–6.
25. Ouzzani M, Hammady H, Fedorowicz Z, Elmagarmid A. Rayyan-a web and mobile app for systematic reviews. *Syst Rev*. 2016;5(1):210.
26. Rahman M, Fukui T. Biomedical research productivity: Factors across the countries. *Int J Technol Assess Health Care*. 2003;19(1):249–52.
27. Hwang TJ, Kesselheim AS, Vokinger KN. Lifecycle Regulation of Artificial Intelligence- And Machine Learning-Based Software Devices in Medicine. Vol. 322, *JAMA - Journal of the American Medical Association*. 2019. p. 2285–6.
28. FDA Center for Devices and Radiological Health. Web-page on Nucleic Acid Based Tests by the Food and Drug Administration (FDA) [Internet]. 2021 [cited 2021 Aug 2]. Available from: <https://www.fda.gov/medical-devices/vitro-diagnostics/nucleic-acid-based-tests>
29. Winner BS, Sgroi DC, Ryan PD, Bruinsma TJ, Glas AM, Male A, et al. Analysis of the mamma print breast cancer assay in a predominantly postmenopausal cohort. *Clin Cancer Res*. 2008;14(10):2988–93.
30. Slodkowska EA, Ross JS. MammaPrint™ 70-gene signature: Another milestone in personalized medical care for breast cancer patients. *Expert Rev Mol Diagn*. 2009;9(5):417–22.
31. Sapino A, Roepman P, Linn SC, Snel MHJ, Delahaye LJMJ, Van Den Akker J, et al. MammaPrint molecular diagnostics on formalin-fixed, paraffin-embedded tissue. *J Mol Diagnostics*. 2014;16(2):190–7.
32. Mook S, Knauer M, Bueno-De-Mesquita JM, Retel VP, Wesseling J, Linn SC, et al. Metastatic potential of T1 breast cancer can be predicted by the 70-gene MammaPrint signature. *Ann Surg Oncol*. 2010;17(5):1406–13.
33. Maak M, Simon I, Nitsche U, Roepman P, Snel M, Glas AM, et al. Independent validation of a prognostic genomic signature (ColoPrint) for patients with stage II colon cancer. *Ann Surg*.

- 2013;257(6):1053–8.
34. Kopetz S, Tabernero J, Rosenberg R, Jiang Z, Moreno V, Bachleitner-Hofmann T, et al. Genomic Classifier ColoPrint Predicts Recurrence in Stage II Colorectal Cancer Patients More Accurately Than Clinical Factors. *Oncologist*. 2015;20(2):127–33.
 35. Tan IB, Tan P. Genetics: An 18-gene signature (ColoPrint®) for colon cancer prognosis. *Nat Rev Clin Oncol*. 2011;8(3):131–3.
 36. Rosenberg R, Maak M, Simon I, Nitsche U, Schuster T, Kuenzli B, et al. Independent validation of a prognostic genomic profile (ColoPrint) for stage II colon cancer (CC) patients. *J Clin Oncol*. 2011;29(4_suppl):358–358.
 37. Salazar R, de Waard JW, Glimelius B, Marshall J, Klaase J, Van Der Hoeven J, et al. The PARSC trial, a prospective study for the assessment of recurrence risk in stage II colon cancer (CC) patients using ColoPrint. *J Clin Oncol*. 2012;30(4_suppl):678–678.
 38. Tabernero J, Moreno V, Rosenberg R, Nitsche U, Bachleitner-Hofmann T, Lanza G, et al. Clinical and technical validation of a genomic classifier (ColoPrint) for predicting outcome of patients with stage II colon cancer. *J Clin Oncol*. 2012;30(4_suppl):384–384.
 39. Bachleitner-Hofmann T, Simon I, Salazar R, Tabernero J, Rosenberg R, van der Akker J, et al. Development and Validation of a Robust Molecular Diagnostic Test (COLOPRINT) for Predicting Outcome in Stage II Colon Cancer Patients. *Ann Oncol*. 2012;23:ix179.
 40. Nielsen T, Wallden B, Schaper C, Ferree S, Liu S, Gao D, et al. Analytical validation of the PAM50-based Prosigna Breast Cancer Prognostic Gene Signature Assay and nCounter Analysis System using formalin-fixed paraffin-embedded breast tumor specimens. *BMC Cancer*. 2014;14(1).
 41. Wallden B, Storhoff J, Nielsen T, Dowidar N, Schaper C, Ferree S, et al. Development and verification of the PAM50-based Prosigna breast cancer gene signature assay. *BMC Med Genomics*. 2015;8(1).
 42. Alvarado MD, Prasad C, Rothney M, Cherbavaz DB, Sing AP, Baehner FL, et al. A Prospective Comparison of the 21-Gene Recurrence Score and the PAM50-Based Prosigna in Estrogen Receptor-Positive Early-Stage Breast Cancer. *Adv Ther*. 2015;32(12):1237–47.
 43. Jensen MB, Lænkholm AV, Nielsen TO, Eriksen JO, Wehn P, Hood T, et al. The Prosigna gene expression assay and responsiveness to adjuvant cyclophosphamide-based chemotherapy in premenopausal high-risk patients with breast cancer. *Breast Cancer Res*. 2018;20(1).
 44. Hequet D, Callens C, Gentien D, Albaud B, Mouret-Reynier MA, Dubot C, et al. Prospective, multicenter French study evaluating the clinical impact of the Breast Cancer Intrinsic Subtype-Prosigna® Test in the management of early-stage breast cancers. *PLoS One*. 2017;12(10).
 45. Parker JS, Mullins M, Cheang MCU, Leung S, Voduc D, Vickery T, et al. Supervised risk predictor of breast cancer based on intrinsic subtypes. *J Clin Oncol*. 2009;27(8):1160–7.
 46. Kelly CM, Krishnamurthy S, Bianchini G, Litton JK, Gonzalez-Angulo AM, Hortobagyi GN, et al. Utility of oncotype DX risk estimates in clinically intermediate risk hormone receptor-positive, HER2-normal, grade II, lymph node-negative breast cancers. *Cancer*. 2010;116(22):5161–7.
 47. Lo SS, Mumby PB, Norton J, Rychlik K, Smerage J, Kash J, et al. Prospective multicenter study of the impact of the 21-gene recurrence score assay on medical oncologist and patient adjuvant breast cancer treatment selection. *J Clin Oncol*. 2010;28(10):1671–6.
 48. Carlson JJ, Roth JA. The impact of the Oncotype Dx breast cancer assay in clinical practice:

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

A systematic review and meta-analysis. Vol. 141, Breast Cancer Research and Treatment. 2013. p. 13–22.

49. Thakur SS, Li H, Chan AMY, Tudor R, Bigras G, Morris D, et al. The use of automated Ki67 analysis to predict Oncotype DX risk-of-recurrence categories in early-stage breast cancer. *PLoS One*. 2018/01/06. 2018;13(1):e0188983.
50. Pease AM, Riba LA, Gruner RA, Tung NM, James TA. Oncotype DX® Recurrence Score as a Predictor of Response to Neoadjuvant Chemotherapy. *Ann Surg Oncol*. 2019;
51. Gianni L, Zambetti M, Clark K, Baker J, Cronin M, Wu J, et al. Gene expression profiles in paraffin-embedded core biopsy tissue predict response to chemotherapy in women with locally advanced breast cancer. *J Clin Oncol*. 2005;23(29):7265–77.
52. Paik S, Shak S, Tang G, Kim C, Baker J, Cronin M, et al. A multigene assay to predict recurrence of tamoxifen-treated, node-negative breast cancer. *N Engl J Med*. 2004;351(27):2817–26.
53. Dowsett M, Sestak I, Lopez-Knowles E, Sidhu K, Dunbier AK, Cowens JW, et al. Comparison of PAM50 risk of recurrence score with oncotype DX and IHC4 for predicting risk of distant recurrence after endocrine therapy. *J Clin Oncol*. 2013;31(22):2783–90.
54. Marrone M, Potosky AL, Penson D, Freedman AN. A 22 gene-expression assay, decipher® (GenomeDx biosciences) to predict five-year risk of metastatic prostate cancer in men treated with radical prostatectomy. *PLoS Curr*. 2015;7(EVIDENCEONGENOMICTESTS).
55. Nguyen PL, Haddad Z, Lam LLC, Ong K, Buerki C, Deheshi S, et al. Evaluation of the Decipher prostate cancer classifier to predict metastasis and disease-specific mortality from genomic analysis of diagnostic prostate needle biopsy specimens. *J Clin Oncol*. 2017;35(6_suppl):4–4.
56. Magi-Galluzzi C, Yousefi K, Haddad Z, Palmer-Aronsten B, Lam LLC, Buerki C, et al. Validation of the Decipher prostate cancer classifier for predicting 10-year postoperative metastasis from analysis of diagnostic needle biopsy specimens. *J Clin Oncol*. 2016;34(2_suppl):59–59.
57. Dalela D, Löppenber B, Sood A, Sammon J, Abdollah F. Contemporary Role of the Decipher® Test in Prostate Cancer Management: Current Practice and Future Perspectives. *Rev Urol*. 2016;18(1):1–9.
58. Klein EA, Haddad Z, Yousefi K, Lam LLC, Wang Q, Choerung V, et al. Decipher Genomic Classifier Measured on Prostate Biopsy Predicts Metastasis Risk. *Urology*. 2016;90:148–52.
59. Weiss LM, Chu P, Schroeder BE, Singh V, Zhang Y, Erlander MG, et al. Blinded comparator study of immunohistochemical analysis versus a 92-gene cancer classifier in the diagnosis of the primary site in metastatic tumors. *J Mol Diagnostics*. 2013;15(2):263–9.
60. Greco FA, Spigel DR, Yardley DA, Erlander MG, Ma X, Hainsworth JD. Molecular Profiling in Unknown Primary Cancer: Accuracy of Tissue of Origin Prediction. *Oncologist*. 2010;15(5):500–6.
61. Hainsworth JD, Rubin MS, Spigel DR, Boccia R V., Raby S, Quinn R, et al. Molecular gene expression profiling to predict the tissue of origin and direct site-specific therapy in patients with carcinoma of unknown primary site: A prospective trial of the Sarah cannon research institute. *J Clin Oncol*. 2013;31(2):217–23.
62. Harrison G, Sosa JA, Jiang X. Evaluation of the Afirma gene expression classifier in repeat indeterminate thyroid nodules. *Arch Pathol Lab Med*. 2017;141(7):985–9.
63. Chudova D, Wilde JI, Wang ET, Wang H, Rabbee N, Egidio CM, et al. Molecular classification of thyroid nodules using high-dimensionality genomic data. *J Clin Endocrinol*

- 1
2
3 Metab. 2010;95(12):5296–304.
4
5 64. Kim MI, Alexander EK. Diagnostic use of molecular markers in the evaluation of thyroid
6 nodules. Vol. 18, *Endocrine Practice*. 2012. p. 796–802.
7
8 65. Ali SZ, Fish SA, Lanman R, Randolph GW, Sosa JA. Use of the Afirma® gene expression
9 classifier for preoperative identification of benign thyroid nodules with indeterminate fine
10 needle aspiration cytopathology. *PLoS Currents*. 2013. p. 1–7.
11
12 66. McIver B, Castro MR, Morris JC, Bernet V, Smallridge R, Henry M, et al. An independent
13 study of a gene expression classifier (Afirma) in the evaluation of cytologically indeterminate
14 thyroid nodules. *J Clin Endocrinol Metab*. 2014;99(11):4069–77.
15
16 67. Lastra RR, Pramick MR, Crammer CJ, LiVolsi VA, Baloch ZW. Implications of a suspicious
17 Afirma test result in thyroid fine-needle aspiration cytology: An institutional experience.
18 *Cancer Cytopathol*. 2014;122(10):737–44.
19
20 68. Kim JY, Park SC, Lee JK, Choi SJ, Hahm KS, Park Y. Novel Antibacterial Activity of β 2-
21 Microglobulin in Human Amniotic Fluid. *PLoS One*. 2012;
22
23 69. Frampton GM, Fichtenholtz A, Otto GA, Wang K, Downing SR, He J, et al. Development
24 and validation of a clinical cancer genomic profiling test based on massively parallel DNA
25 sequencing. *Nat Biotechnol*. 2013;31(11):1023–31.
26
27 70. Wang K, Sanchez-Martin M, Wang X, Knapp KM, Koche R, Vu L, et al. Patient-derived
28 xenotransplants can recapitulate the genetic driver landscape of acute leukemias.
29 *Leukemia*. 2017;31(1):151–8.
30
31 71. Tarlock K, He J, Zhong S, Ries RE, Bailey M, Morley S, et al. Distinct age-associated
32 genomic profiles in acute myeloid leukemia (AML) using FoundationOne heme. *J Clin
33 Oncol*. 2016;34(15_suppl):7041–7041.
34
35 72. Lieber DS, Kennedy MR, Johnson DB, Greenbowe JR, Frampton GM, Schrock AB, et al.
36 Abstract B16: Validation and clinical feasibility of a Foundation Medicine assay to identify
37 immunotherapy response potential through tumor mutational burden (TMB). In 2017.
38
39 73. Lee Deak K, Jackson JB, Valkenburg KC, Keefer LA, Robinson Gerding KM, Angiuoli SV, et
40 al. Next-Generation Sequencing Concordance Analysis of Comprehensive Solid Tumor
41 Profiling between a Centralized Specialty Laboratory and the Decentralized PGDx Elio
42 Tissue Complete Kitted solution. *J Mol Diagnostics [Internet]*. 2021 Jul;in press. Available
43 from: <https://linkinghub.elsevier.com/retrieve/pii/S1525157821002105>
44
45 74. Labriola MK, Zhu J, Gupta R, McCall S, Jackson J, Kong EF, et al. Characterization of
46 tumor mutation burden, PD-L1 and DNA repair genes to assess relationship to immune
47 checkpoint inhibitors response in metastatic renal cell carcinoma. *J Immunother Cancer
48 [Internet]*. 2020 Mar;8(1):e000319. Available from:
49 <https://jitc.bmj.com/lookup/doi/10.1136/jitc-2019-000319>
50
51 75. Yamani MH, Taylor DO, Rodriguez ER, Cook DJ, Zhou L, Smedira N, et al. Transplant
52 Vasculopathy Is Associated With Increased AlloMap Gene Expression Score. *J Hear Lung
53 Transplant*. 2007;26(4):403–6.
54
55 76. Yamani MH, Taylor DO, Haire C, Smedira N, Starling RC. Post-transplant ischemic injury is
56 associated with up-regulated AlloMap gene expression. *Clin Transplant*. 2007;21(4):523–5.
57
58 77. Kobashigawa J, Patel J, Azarbal B, Kittleson M, Chang D, Czer L, et al. Randomized Pilot
59 Trial of Gene Expression Profiling Versus Heart Biopsy in the First Year after Heart
60 Transplant: Early Invasive Monitoring Attenuation Through Gene Expression Trial. *Circ Hear
Fail*. 2015;8(3):557–64.
78. Wingrove JA, Daniels SE, Sehnert AJ, Tingley W, Elashoff MR, Rosenberg S, et al.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Correlation of peripheral-blood gene expression with the extent of coronary artery stenosis. *Circ Cardiovasc Genet*. 2008;1(1):31–8.

79. Rosenberg S, Dehais C, Ducray F, Alentron A, Tanguy M, De Reyneis A, et al. OS11.3 Machine learning for better prognostic stratification and driver genes identification in 1p/19q-codeleted grade III gliomas. *Neuro Oncol* [Internet]. 2017;19(suppl_3):iii22–iii22. Available from: <http://dx.doi.org/10.1093/neuonc/nox036>
80. Vargas J, Lima JAC, Kraus WE, Douglas PS, Rosenberg S. Use of the Corus® CAD Gene Expression Test for Assessment of Obstructive Coronary Artery Disease Likelihood in Symptomatic Non-Diabetic Patients. *PLoS Currents*. 2013.
81. Elashoff MR, Wingrove JA, Beineke P, Daniels SE, Tingley WG, Rosenberg S, et al. Development of a blood-based gene expression algorithm for assessment of obstructive coronary artery disease in non-diabetic patients. *BMC Med Genomics*. 2011;4.
82. Rosenberg S, Elashoff MR, Beineke P, Daniels SE, Wingrove JA, Tingley WG, et al. Multicenter validation of the diagnostic accuracy of a blood-based gene expression test for assessing obstructive coronary artery disease in nondiabetic patients. *Ann Intern Med*. 2010;153(7):425–34.
83. Brahe CH, Østergaard M, Johansen JS, Defranoux N, Wang X, Bolce R, et al. Predictive value of a multi-biomarker disease activity score for clinical remission and radiographic progression in patients with early rheumatoid arthritis: a post-hoc study of the OPERA trial. *Scand J Rheumatol*. 2019;
84. Chernoff D, Scott Eastman P, Hwang CC, Flake DD, Wang X, Kivitz A, et al. Determination of the minimally important difference (MID) in multi-biomarker disease activity (MBDA) test scores: impact of diurnal and daily biomarker variation patterns on MBDA scores. *Clin Rheumatol*. 2019;38(2):437–45.
85. Curtis JR, Weinblatt ME, Shadick NA, Brahe CH, Østergaard M, Hetland ML, et al. Validation of the adjusted multi-biomarker disease activity score as a prognostic test for radiographic progression in rheumatoid arthritis: a combined analysis of multiple studies. *Arthritis Res Ther*. 2021;
86. Curtis JR, Xie F, Yang S, Danila MI, Owensby JK, Chen L. Uptake and Clinical Utility of Multibiomarker Disease Activity Testing in the United States. *J Rheumatol* [Internet]. 2019;46(3):237–44. Available from: <http://www.ncbi.nlm.nih.gov/pubmed/30442830>
87. Curtis JR, Van Der Helm-Van Mil AH, Knevel R, Huizinga TW, Haney DJ, Shen Y, et al. Validation of a novel multibiomarker test to assess rheumatoid arthritis disease activity. *Arthritis Care Res*. 2012;64(12):1794–803.
88. Food and Drug Administration. Helix Genetic Health Risk App For Late-Onset Alzheimer's Disease - FDA Review Decision Summary. 2020; Available from: <https://www.accessdata.fda.gov/scripts/cdrh/cfdocs/cfpmn/pmn.cfm?ID=K192073>
89. Cirulli ET, White S, Read RW, Elhanan G, Metcalf WJ, Tanudjaja F, et al. Genome-wide rare variant analysis for thousands of phenotypes in over 70,000 exomes from two cohorts. *Nat Commun*. 2020;
90. Lu JT, Ferber M, Hagenkord J, Levin E, South S, Kang HP, et al. Evaluation for Genetic Disorders in the Absence of a Clinical Indication for Testing: Elective Genomic Testing. *Journal of Molecular Diagnostics*. 2019.
91. Grzymalski JJ, Elhanan G, Morales Rosado JA, Smith E, Schlauch KA, Read R, et al. Population genetic screening efficiently identifies carriers of autosomal dominant diseases. *Nat Med*. 2020;
92. Lemke C, Budka M, Gabrys B. Metalearning: a survey of trends and technologies. *Artif Intell*

1
2
3 Rev. 2015;
4

- 5 93. Sima C, Braga-Neto UM, Dougherty ER. High-dimensional bolstered error estimation.
6 *Bioinformatics*. 2011;
7
8 94. Ali M, Aittokallio T. Machine learning and feature selection for drug response prediction in
9 precision oncology applications. *Biophysical Reviews*. 2019.
10
11 95. Barla A, Jurman G, Riccadonna S, Merler S, Chierici M, Furlanello C. Machine learning
12 methods for predictive proteomics. *Br Bioinform*. 2008/03/04. 2008;9(2):119–28.
13
14 96. Grollemund V, Pradat PF, Querin G, Delbot F, Le Chat G, Pradat-Peyre JF, et al. Machine
15 learning in amyotrophic lateral sclerosis: Achievements, pitfalls, and future directions. *Front*
16 *Neurosci* [Internet]. 2019;13. Available from: <http://dx.doi.org/10.3389/fnins.2019.00135>
17
18 97. Ho DSW, Schierding W, Wake M, Saffery R, O'Sullivan J. Machine learning SNP based
19 prediction for precision medicine. *Front Genet* [Internet]. 2019;10(MAR). Available from:
20 <http://dx.doi.org/10.3389/fgene.2019.00267>
21
22 98. Krittanawong C, Bomback AS, Baber U, Bangalore S, Messerli FH, Wilson Tang WH. Future
23 Direction for Using Artificial Intelligence to Predict and Manage Hypertension. *Curr*
24 *Hypertens Rep*. 2018/07/08. 2018;20(9):75.
25
26 99. Long NP, Yoon SJ, Anh NH, Nghi TD, Lim DK, Hong YJ, et al. A systematic review on
27 metabolomics-based diagnostic biomarker discovery and validation in pancreatic cancer.
28 *Metabolomics*. 2019/03/05. 2018;14(8):109.
29
30 100. Martinez BI, Stabenfeldt SE. Current trends in biomarker discovery and analysis tools for
31 traumatic brain injury. *J Biol Eng* [Internet]. 2019;13(1). Available from:
32 <http://dx.doi.org/10.1186/s13036-019-0145-8>
33
34 101. Patil S, Awan KH, Arakeri G, Seneviratne CJ, Muddur N, Malik S, et al. Machine learning
35 and its potential applications to the genomic study of head and neck cancer-A systematic
36 review. *J Oral Pathol Med*. 2019;48(9):773–9.
37
38 102. Saini G, Mittal K, Rida P, Janssen EAM, Gogineni K, Aneja R. Panoptic view of prognostic
39 models for personalized breast cancer management. *Cancers (Basel)* [Internet]. 2019;11(9).
40 Available from: <http://dx.doi.org/10.3390/cancers11091325>
41
42 103. European Council. Council conclusions on personalised medicine for patients. *Off J Eur*
43 *Union* [Internet]. 2015;431(2):1–4. Available from: [https://eur-lex.europa.eu/legal-](https://eur-lex.europa.eu/legal-content/EN/TXT/PDF/?uri=CELEX:52015XG1217(01)&from=EN)
44 [content/EN/TXT/PDF/?uri=CELEX:52015XG1217\(01\)&from=EN](https://eur-lex.europa.eu/legal-content/EN/TXT/PDF/?uri=CELEX:52015XG1217(01)&from=EN)
45
46 104. Labriola MK, Zhu J, Gupta R, McCall S, Jackson J, Kong EF, et al. Characterization of
47 tumor mutation burden, PD-L1 and DNA repair genes to assess relationship to immune
48 checkpoint inhibitors response in metastatic renal cell carcinoma. *J Immunother Cancer*.
49 2020;8(1).
50
51
52
53
54
55
56
57
58
59
60

Author Contributions

Study conception and design: EG, AR.

Methodology: CG, RB.

Data collection and analysis: EG, AR

Original draft preparation: EG

Review and editing: AR, EG, PG, CG, JDM, RB.

Project supervision: PG

Funding acquisition: JDM.

All authors have read and revised the manuscript and approved the final version.

The members of the PERMIT group were involved in the preparation or revision of the joint protocol of the four scoping reviews of the PERMIT series, attended the joint workshop (consultation exercise) and are co-authors of the other scoping reviews of the PERMIT series.

Collaborators

PERMIT group:

1. Antonio L. Andreu
2. Florence Bietrix,
3. Florie Brion Bouvier
4. Montserrat Carmona Rodriguez
5. Maria del Mar Polo-de Santos,
6. Maddalena Fratelli,
7. Rainer Girgenrath,
8. Alexander Grundmann,
9. Josep Maria Haro,
10. Frank Hulstaert,
11. Iñaki Imaz-Iglesia,
12. Setefilla Luengo Matos
13. Emmet McCormack,
14. Albert Sanchez Niubo,
15. Emanuela Oldoni,
16. Raphael Porcher,
17. Vibeke Fosse,
18. Luis M. Sánchez-Gómez,
19. Lorena San Miguel,
20. Cecilia Superchi,
21. Teresa Torres,
22. Anna Monistrol Mula

Funding statement

This project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No. 874825.

Competing interests statement

None declared

Ethics approval

This study was based entirely on a scoping review of relevant published literature and did not require an ethics approval.

Patient consent

This study did not require consent from patients, because it does not use individual data.

Permission to reproduce material from other sources

This study has cited all references which are published and publicly available.

Data sharing statement

The study protocol was published on the online platform Zenodo (19). Copies of searches and data extraction sheets will be made publicly available on Zenodo as part of the database collection for all scoping reviews conducted in the PERMIT project.

Figure legends

Fig. 1. Keyword based search strategy for the scoping review. Four categories of keywords were defined to retrieve relevant articles from the biomedical literature on machine learning analyses of omics data for personalised medicine, which include a validation study (highlighted by the coloured boxes in the centre). For each category relevant keywords were determined, including controlled vocabulary terms from the Medical Subject Headings (MeSH) thesaurus by the US National Library of Medicine (upper and lower boxes with frames coloured according to the corresponding category). As indicated by the keyword “AND” in the centre, a conjunctive search was conducted, i.e., every retrieved article had to contain at least one keyword from each category. This strategy was adapted for searching the other databases.

Fig. 2. Study selection flow diagram. Flow diagram of the procedure for the scoping review article identification, screening, eligibility assessment, and final inclusion, according to the PRISMA scheme (23). Reasons for excluding full-text were not mutually exclusive.

Fig. 3. Validation methods used in omics biomarker studies. Stacked bar chart of the number of articles retrieved in the scoping review for different categories of validation methods used in the underlying biomarker studies (covering time periods from 2000 to 2021). The majority of studies use only internal cohort validation approaches, such as cross-validation (CV), training/test set split validation, resampling/bootstrapping-based validation, out-of-bag validation (for tree-based classifiers), and combinations of CV and test set validation within the same cohort. Studies with an external validation on an independent patient cohort (with or without an additional internal cross-validation) are still underrepresented, even in more recent time periods. All filtered full-text articles derived from the scoping review except for review articles were included in the analysis.

Fig. 4. Map representation of country statistics for the selected articles. The number of articles originating from different countries among the studies selected in the full-text review are visualized on a world map representation using a colour gradient from blue (1 article) to red (98 articles = maximum contribution by a single country; using a logarithmic colour gradient scale to highlight differences over a broad value range).

Fig. 5. Representation of study types among the selected articles. The percentage of articles describing case-control studies, therapy/drug response studies, differential diagnosis studies, prognostic and survival prediction studies, as well as review studies and other study types is represented as a pie chart.

Fig. 6. Characteristics of successful omics-based studies. Six main categories of design and implementation aspects that characterize successful omics-based biomarker development studies were identified (starting from the centre left in the figure and proceeding clockwise): 1) Adequacy of the study design & sample size selection; 2) Rigor and robustness of the statistical evaluation; 3) Clarity of scope and goals; 4) Completeness and reproducibility of the study documentation; 5) Interpretability and biological plausibility of the created predictive models; 6) Integration of prior biological knowledge into the model building and validation procedures.

Tables

Name	Test approval type	Purpose	References
MammaPrint®	FDA-cleared Assay	breast cancer risk-of-recurrence assessment	(6,29–32)
ColoPrint®	LDT	colon cancer development of distant metastasis prediction	(33–38)
Prosigna® Assay / PAM50	FDA-cleared Assay	breast cancer risk of distant recurrence prediction	(40–44)
Oncotype DX®	LDT	breast cancer risk-of-recurrence assessment	(8,46–49)
Decipher®	LDT	prostate cancer metastatic risk prediction	(9,54–58)
Cancer Type ID®	LDT	predict tumour type for cancers of unknown / uncertain diagnosis	(15,59–61)
Afirma™ Gene Expression Classifier	LDT	discriminate between benign and cancerous thyroid nodules	(62–67)
Foundation One® Heme	LDT	test for hematologic malignancies, sarcomas, or solid tumours	(14,69–71)
PGDx Elio™ Tissue Complete	FDA-cleared Assay	test to assess somatic mutations and tumour mutation burden for solid tumours	(73,104)
AlloMap® Heart	FDA-cleared Assay	identifying heart transplant recipients with risk of cellular rejection	(13,75–77)
Corus® CAD	LDT	identify obstructive coronary artery disease	(11,78–81)
Vectra® DA	LDT	multi-biomarker blood test for rheumatoid arthritis	(83–86)
Helix® Laboratory Platform & Health Risk App for Late-onset Alzheimer's	FDA-cleared medical device	whole exome sequencing constituent device based for reporting and interpreting general genetic health risks	(89–91)

Tab. 1. Examples of clinically approved omics-derived diagnostic or prognostic tests designs applied to personalised medicine (synonyms for the same test are separated by the “/”-symbol). FDA-approval status was checked on the web-site by the FDA (28) and reflects the status as of July 2021.

Page 29 of 40
stratified medicine" OR biomarker* OR "precision medicine"
OR "personalized medicine" OR "personalised medicine"
OR "individualized Medicine" OR "individualised Medicine"
OR "individualized therapy" OR "individualised therapy" OR
"patient stratification" OR pharmacogenetics OR "patient
specific modeling" OR "personalized clinical decision
making" OR "personalised clinical decision making" OR
"prediction of response" OR "prediction of responses" OR
"Biomarkers"[Mesh] OR "Precision Medicine"[Mesh]

BMJ Open
Genomics"[Mesh]) OR "Metabolomics"[Mesh]) OR
"Epigenomics"[Mesh]) OR "Microarray Analysis"[Mesh]) OR
"Mass Spectrometry"[Mesh] OR Omic* OR "omic based"
OR "multi omic" OR "multi omics" OR genomic* OR
transcriptomic* OR proteomic* OR metabolomic OR
lipidomic* OR epigenomic* OR microarray OR "RNA seq"
OR "mass spectrometry")

PERSONALISED
MEDICINE

OMICS

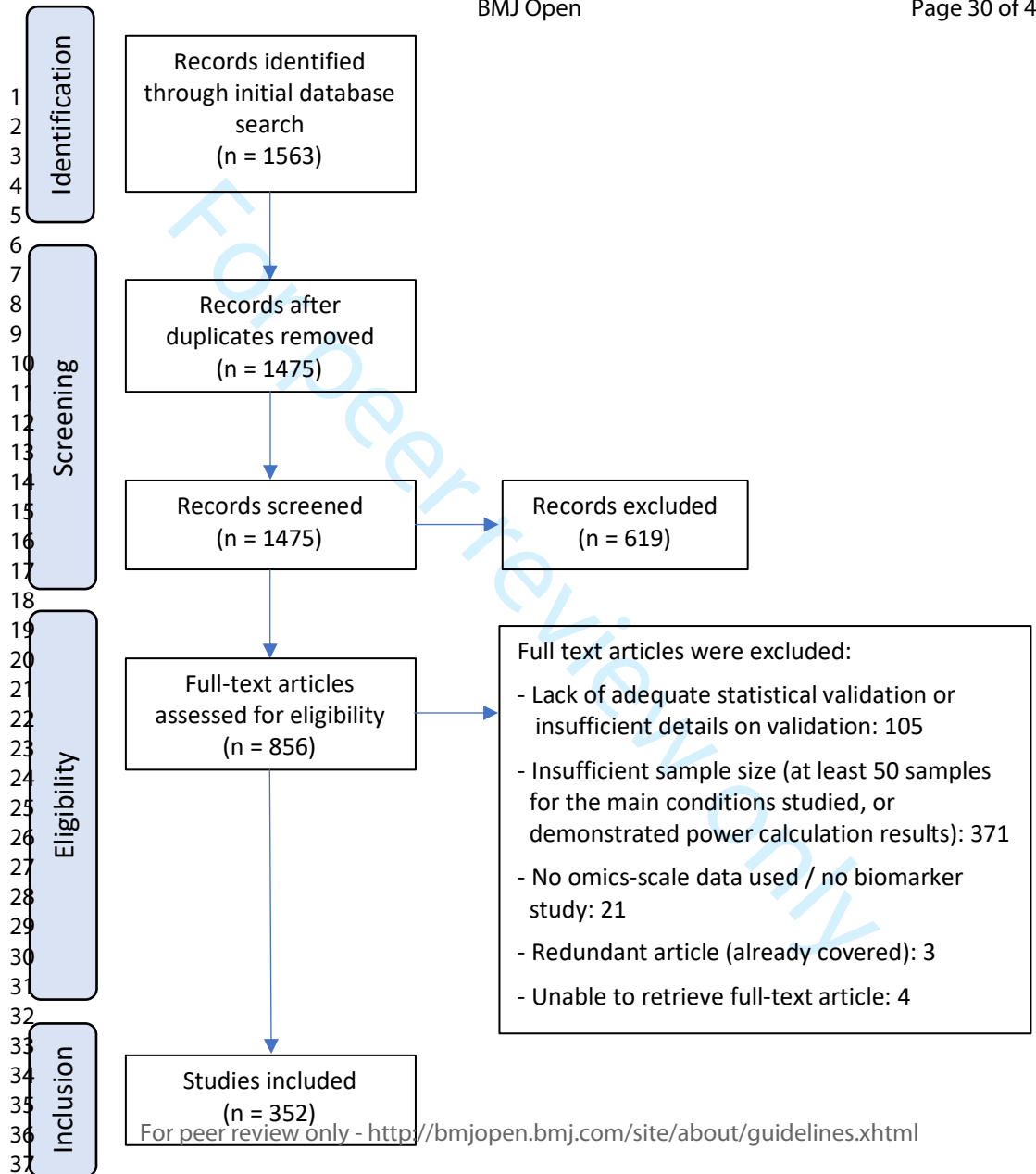
AND

MACHINE
LEARNING

VALIDATION

Machine Learning"[Mesh] OR «Machine learning" OR
"statistical learning" OR "supervised learning" OR
"unsupervised learning"

Validation Studies as Topic"[Mesh]) OR "Validation Study"
[Publication Type] OR "Sensitivity and Specificity"[Mesh])
OR "Benchmarking"[Mesh]) OR validation OR validity OR
validated OR "cross validation" "cross validated" OR
"clinical utility*" OR accuracy OR robustness OR reliability*
OR sensitivity OR specificity OR benchmark* OR bias OR
"cross study" OR "cross studies")



Page 31 of 40
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38

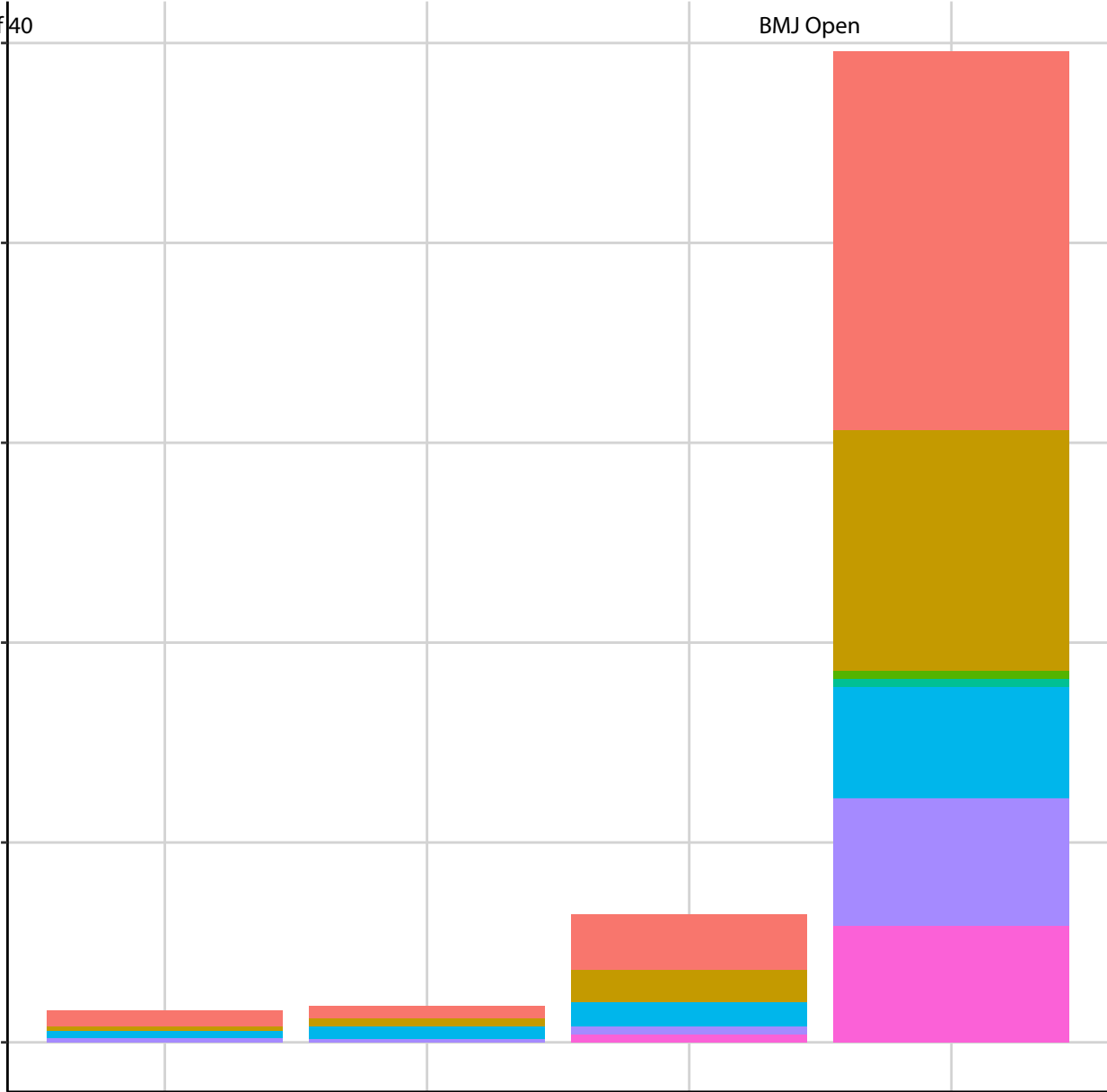
Validation methods

Internal validation:

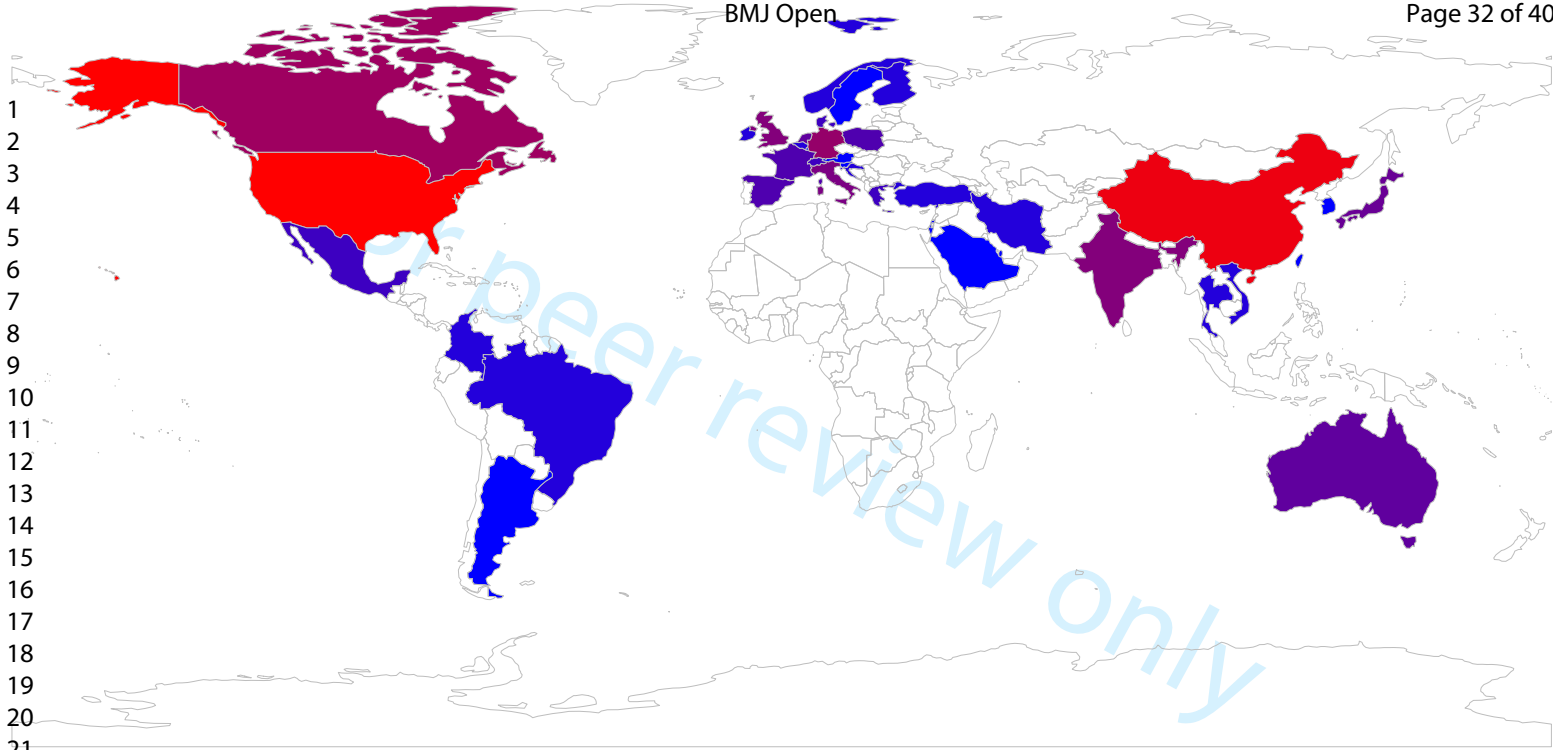
- Cross-validation (CV)
- Training/test set validation
- Resampling
- Out-of-bag internal validation
- CV + internal cohort validation

External validation:

- CV + external cohort validation
- External cohort validation



1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26

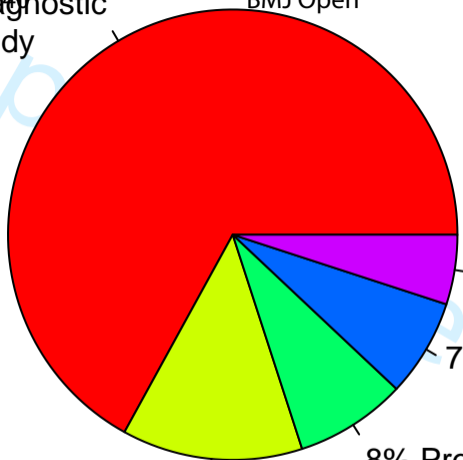


For peer review only <http://bmjopen.bmj.com/site/about/guidelines.xhtml>

Page 3 of 40
67% Diagnostic study

BMJ Open

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16



5% Others

7% Therapy/drug response study

8% Prognostic/survival study

13% Review study

Characteristics of successful omics-derived biomarker studies

Statistical evaluation

- cross-validation & external testing
- adequacy of performance & robustness metrics
- multiple testing correction

Clarity of scope & goals

- Inclusion/exclusion criteria
- primary/secondary outcomes

Study design & sample size

- statistical power
- balanced study groups
- batch effect avoidance/correction
- matching/adjustment for confounders/biases

Study documentation

- instruments, settings & parameters
- reproducible methods description
- data provenance

Integration of prior knowledge

- molecular pathways & networks
- clinical data & real-world data
- multi-omics data

Model interpretability

- explainable predictions & human-interpretable models
- biological plausibility & mechanistic understanding

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29

Online Supplementary file 1 – PRISMA-ScR Checklist

Preferred Reporting Items for Systematic reviews and Meta-Analyses extension for Scoping Reviews (PRISMA-ScR) Checklist (17).

SECTION	ITEM	PRISMA-ScR CHECKLIST ITEM	REPORTED ON PAGE #
TITLE			
Title	1	Identify the report as a scoping review.	1
ABSTRACT			
Structured summary	2	Provide a structured summary that includes (as applicable): background, objectives, eligibility criteria, sources of evidence, charting methods, results, and conclusions that relate to the review questions and objectives.	2
INTRODUCTION			
Rationale	3	Describe the rationale for the review in the context of what is already known. Explain why the review questions/objectives lend themselves to a scoping review approach.	4
Objectives	4	Provide an explicit statement of the questions and objectives being addressed with reference to their key elements (e.g., population or participants, concepts, and context) or other relevant key elements used to conceptualize the review questions and/or objectives.	4
METHODS			
Protocol and registration	5	Indicate whether a review protocol exists; state if and where it can be accessed (e.g., a Web address); and if available, provide registration information, including the registration number.	5
Eligibility criteria	6	Specify characteristics of the sources of evidence used as eligibility criteria (e.g., years considered, language, and publication status), and provide a rationale.	5-6
Information sources*	7	Describe all information sources in the search (e.g., databases with dates of coverage and contact with authors to identify additional sources), as well as the date the most recent search was executed.	5
Search	8	Present the full electronic search strategy for at least 1 database, including any limits used, such that it could be repeated.	5 (Online Suppl. File 2)
Selection of sources of evidence†	9	State the process for selecting sources of evidence (i.e., screening and eligibility) included in the scoping review.	6

SECTION	ITEM	PRISMA-ScR CHECKLIST ITEM	REPORTED ON PAGE #
Data charting process‡	10	Describe the methods of charting data from the included sources of evidence (e.g., calibrated forms or forms that have been tested by the team before their use, and whether data charting was done independently or in duplicate) and any processes for obtaining and confirming data from investigators.	6-7
Data items	11	List and define all variables for which data were sought and any assumptions and simplifications made.	6 (Online Suppl. File 3)
Critical appraisal of individual sources of evidence§	12	If done, provide a rationale for conducting a critical appraisal of included sources of evidence; describe the methods used and how this information was used in any data synthesis (if appropriate).	Click here to enter text.
Synthesis of results	13	Describe the methods of handling and summarizing the data that were charted.	6-7
RESULTS			
Selection of sources of evidence	14	Give numbers of sources of evidence screened, assessed for eligibility, and included in the review, with reasons for exclusions at each stage, ideally using a flow diagram.	7 (Fig. 2)
Characteristics of sources of evidence	15	For each source of evidence, present characteristics for which data were charted and provide the citations.	7 (Table 1)
Critical appraisal within sources of evidence	16	If done, present data on critical appraisal of included sources of evidence (see item 12).	Click here to enter text.
Results of individual sources of evidence	17	For each included source of evidence, present the relevant data that were charted that relate to the review questions and objectives.	7-13
Synthesis of results	18	Summarize and/or present the charting results as they relate to the review questions and objectives.	7-13
DISCUSSION			
Summary of evidence	19	Summarize the main results (including an overview of concepts, themes, and types of evidence available), link to the review questions and objectives, and consider the relevance to key groups.	13
Limitations	20	Discuss the limitations of the scoping review process.	13-14
Conclusions	21	Provide a general interpretation of the results with respect to the review questions and	14-15

SECTION	ITEM	PRISMA-ScR CHECKLIST ITEM	REPORTED ON PAGE #
		objectives, as well as potential implications and/or next steps.	
FUNDING			
Funding	22	Describe sources of funding for the included sources of evidence, as well as sources of funding for the scoping review. Describe the role of the funders of the scoping review.	24

For peer review only

Online Supplementary file 2 – Search strategy

Keyword searches conducted in the databases *PubMed*, *EMBASE* and *Web of Science* as part of the scoping review.

1) PubMed Query

Search: (((("Machine learning" OR "statistical learning" OR "supervised learning" OR "unsupervised learning") OR ("Machine Learning"[Mesh])) AND (("Biomarkers"[Mesh]) OR "Precision Medicine"[Mesh])) AND (((Omic* OR "omic based" OR "multi omic" OR "multi omics" OR genomic* OR transcriptomic* OR proteomic* OR metabolomic* OR lipidomic* OR epigenomic* OR microarray OR "RNA seq" OR "mass spectrometry")) OR ("Genomics"[Mesh] OR "Metabolomics"[Mesh] OR "Epigenomics"[Mesh] OR "Microarray Analysis"[Mesh] OR "Mass Spectrometry"[Mesh])) AND ((validation OR validity OR validated OR "cross validation" "cross validated" OR "clinical utility*" OR accuracy OR robustness OR reliability* OR sensitivity OR specificity OR benchmark* OR bias OR "cross study" OR "cross studies")) AND ("2000/01/01"[Date - Entry]: "2021/07/20"[Date - Entry])) Filters: English, French, Italian, Spanish

2) Embase Query

#25: #24 AND [embase]/lim NOT [medline]/lim
 #24: #23 AND [2000-2021]/py
 #23: #20 AND #21 AND ([english]/lim OR [french]/lim OR [italian]/lim OR [spanish]/lim)
 #22: #20 AND #21
 #21: omic*:ti,ab OR 'machine learning':ti,ab OR 'personalized medicine':ti,ab OR 'personalised medicine':ti,ab
 #20: #4 AND #10 AND #16 AND #19
 #19: #17 OR #18
 #18: validation:ti,ab OR validity:ti,ab OR validated:ti,ab OR 'cross validation':ti,ab OR 'cross validated':ti,ab OR test*:ti,ab OR 'clinical utility*':ti,ab OR accuracy:ti,ab OR robustness:ti,ab OR reliability*:ti,ab OR sensitivity:ti,ab OR specificity:ti,ab OR benchmark*:ti,ab OR bias:ti,ab OR 'cross study:ti,ab' OR 'cross studies':ti,ab
 #17: 'validation study'/exp OR 'reliability'/exp OR 'sensitivity and specificity'/exp OR 'benchmarking'/exp
 #16: #14 OR #15
 #15: omic*:ti,ab OR 'omic based':ti,ab OR 'multi omic*':ti,ab OR genomic*:ti,ab OR transcriptomic*:ti,ab OR proteomic*:ti,ab OR metabolomic*:ti,ab OR lipidomic*:ti,ab OR epigenomic*:ti,ab OR microarray:ti,ab OR 'rna seq':ti,ab OR 'mass spectrometr*':ti,ab
 #14: #11 OR #12 OR #13
 #13: 'mass spectrometry'/exp
 #12: 'microarray analysis'/exp
 #11: 'omics'/exp OR 'genomics'/exp OR 'epigenetics'/exp
 #10: #5 OR #6 OR #7 OR #8 OR #9
 #9: 'individualized medicine':ti,ab OR 'individualised medicine':ti,ab OR 'individualized therapy':ti,ab OR 'individualised therapy':ti,ab
 #8: 'personalised medicine':ti,ab
 #7: 'personalized medicine':ti,ab
 #6: 'stratified medicine':ti,ab OR cluster*:ti,ab OR 'sub group*':ti,ab OR subgroup*:ti,ab OR biomarker*:ti,ab OR diagnos*:ti,ab OR prognos*:ti,ab OR 'precision medicine':ti,ab
 #5: 'biological marker'/exp OR 'personalized medicine'/exp
 #4: #1 OR #2 OR #3
 #3: 'machine learning'/exp
 #2: 'statistical learning'/exp

1
2
3 #1: 'machine learning':ti,ab OR 'statistical learning':ti,ab OR 'supervised learning':ti,ab OR
4 'unsupervised learning':ti,ab
5
6

7 **3) Web of Science Query**

8 (((((#1) AND #1) AND #2) AND #3) AND #4) AND #5

9 5: (ALL=(((validation OR validity OR validated OR "cross validation" "cross validated" OR
10 "clinical utility*" OR accuracy OR robustness OR reliability* OR sensitivity OR specificity OR
11 benchmark*

12 OR bias OR "cross study" OR "cross studies")))) AND ALL=(((omic* OR "machine learning"
13 OR

14 "personalized medicine" OR "personalised Medicine"))))

15 4: ALL=(((validation OR validity OR validated OR "cross validation" "cross validated" OR
16 "clinical utility*" OR accuracy OR robustness OR reliability* OR sensitivity OR specificity OR
17 benchmark*

18 OR bias OR "cross study" OR "cross studies"))))

19 3: ALL=(TOPIC: ((Omic* OR "omic based" OR "multi omic*" OR genomic* OR transcriptomic*
20 OR proteomic* OR metabolomic* OR lipidomic* OR epigenomic* OR microarray OR "RNA
21 seq"

22 OR "mass spectrometr*"))))

23 2: (ALL=(TOPIC: (("Machine learning" OR "statistical learning" OR "supervised learning" OR
24 "unsupervised learning")))) AND ALL=(TOPIC: (("stratified medicine" OR cluster* OR "sub

25 group*" OR Subgroup* OR biomarker* OR diagnos* OR prognos* OR "precision medicine"
26 OR

27 "personalized medicine"OR "personalised medicine" OR "individualized Medicine" OR
28 "individualised Medicine" OR "individualized therapy" OR "individualised therapy"))))

29 1: ALL=(TOPIC: (("Machine learning" OR "statistical learning" OR "supervised learning" OR
30 "unsupervised learning"))))
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Online Supplementary file 3 – Data extraction form

Data items extracted from each processed article during the full-text scoping review, and associated qualifications for each item.

Item	Qualifications
Authors	
Title	
Journal	
Volume	
Issue	(if applicable)
Pages	(if applicable)
Year	
Location	
URL / DOI	
Type of publication	<ul style="list-style-type: none"> • Research article • Meeting abstract • Review
Study population and sample size	(if applicable)
Methodology / Study Design	<ul style="list-style-type: none"> • Case-control study • Cases only stratification study <p>(+ further qualification, e.g. treatment response prediction, tumor subtype categorization, recurrence/relapse prediction, survival prediction, tissue-of-origin prediction)</p>
Outcome assessment	<ul style="list-style-type: none"> • Performance measures (e.g. accuracy, sensitivity, specificity, Kohen's Kappa, F-score, AUC) • Validation scheme (cross-validation approach, external validation approach, single cohort or multiple cohorts)
Generic machine learning category	<ul style="list-style-type: none"> • Supervised learning • Unsupervised learning • Other / mixed approaches
Name of specific machine learning approach	(if applicable)

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Main results / key findings that relate to the research question	short description
---	-------------------

For peer review only

BMJ Open

Biomarker discovery studies for patient stratification using machine learning analysis of omics data: a scoping review

Journal:	<i>BMJ Open</i>
Manuscript ID	bmjopen-2021-053674.R2
Article Type:	Original research
Date Submitted by the Author:	06-Nov-2021
Complete List of Authors:	Glaab, Enrico; University of Luxembourg, Luxembourg Centre for Systems Biomedicine Rauschenberger, Armin; University of Luxembourg, Luxembourg Centre for Systems Biomedicine Banzi, Rita; Mario Negri Institute for Pharmacological Research, Center for Health Regulatory Policies Gerardi, Chiara; Mario Negri Institute for Pharmacological Research, Center for Health Regulatory Policies Garcia, Paula; ECRIN, European Clinical Research Infrastructure Network Demotes, Jacques; ECRIN, European Clinical Research Infrastructure Network
Primary Subject Heading:	Patient-centred medicine
Secondary Subject Heading:	Diagnostics, Research methods
Keywords:	BIOTECHNOLOGY & BIOINFORMATICS, NATURAL SCIENCE DISCIPLINES, STATISTICS & RESEARCH METHODS

SCHOLARONE™
Manuscripts



I, the Submitting Author has the right to grant and does grant on behalf of all authors of the Work (as defined in the below author licence), an exclusive licence and/or a non-exclusive licence for contributions from authors who are: i) UK Crown employees; ii) where BMJ has agreed a CC-BY licence shall apply, and/or iii) in accordance with the terms applicable for US Federal Government officers or employees acting as part of their official duties; on a worldwide, perpetual, irrevocable, royalty-free basis to BMJ Publishing Group Ltd ("BMJ") its licensees and where the relevant Journal is co-owned by BMJ to the co-owners of the Journal, to publish the Work in this journal and any other BMJ products and to exploit all rights, as set out in our [licence](#).

The Submitting Author accepts and understands that any supply made under these terms is made by BMJ to the Submitting Author unless you are acting as an employee on behalf of your employer or a postgraduate student of an affiliated institution which is paying any applicable article publishing charge ("APC") for Open Access articles. Where the Submitting Author wishes to make the Work available on an Open Access basis (and intends to pay the relevant APC), the terms of reuse of such Open Access shall be governed by a Creative Commons licence – details of these licences and which [Creative Commons](#) licence will apply to this Work are set out in our licence referred to above.

Other than as permitted in any relevant BMJ Author's Self Archiving Policies, I confirm this Work has not been accepted for publication elsewhere, is not being considered for publication elsewhere and does not duplicate material already published. I confirm all authors consent to publication of this Work and authorise the granting of this licence.

Title

Biomarker discovery studies for patient stratification using machine learning analysis of omics data: a scoping review

Authors

Enrico Glaab^{1,*}, Armin Rauschenberger¹, Rita Banzi², Chiara Gerardi², Paula Garcia³, Jacques Demotes-Mainard³, and the PERMIT Group

Affiliations

¹Luxembourg Centre for Systems Biomedicine, University of Luxembourg, Campus Belval, Esch-sur-Alzette, Luxembourg.

²Center for Health Regulatory Policies, Istituto di Ricerche Farmacologiche Mario Negri, Milan, Italy.

³European Clinical Research Infrastructure Network (ECRIN), Paris, France

*Correspondence: enrico.glaab@uni.lu; 14, avenue du Rock'n'Roll, L-4361 Esch-sur-Alzette, Luxembourg; Phone: +352-466644 6186

Word count

Abstract: 300

Main text: 7153

Keywords

Biomarkers, Scoping Review, Omics, Machine Learning, Stratification

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Abstract

Objective: To review biomarker discovery studies using omics data for patient stratification which led to clinically validated FDA-cleared tests or laboratory developed tests, in order to identify common characteristics and derive recommendations for future biomarker projects.

Design: Scoping review.

Methods: We searched PubMed, EMBASE and Web of Science to obtain a comprehensive list of articles from the biomedical literature published between January 2000 to July 2021, describing clinically validated biomarker signatures for patient stratification, derived using statistical learning approaches. All documents were screened to retain only peer-reviewed research articles, review articles, or opinion articles, covering supervised and unsupervised machine learning applications for omics-based patient stratification. Two reviewers independently confirmed the eligibility. Disagreements were solved by consensus. We focused the final analysis on omics-based biomarkers which achieved the highest level of validation, i.e., clinical approval of the developed molecular signature as a laboratory developed test or FDA approved tests.

Results: Overall, 352 articles fulfilled the eligibility criteria. The analysis of validated biomarker signatures identified multiple common methodological and practical features that may explain the successful test development and guide future biomarker projects. These include study design choices to ensure sufficient statistical power for model building and external testing, suitable combinations of non-targeted and targeted measurement technologies, the integration of prior biological knowledge, strict filtering and inclusion/exclusion criteria, and the adequacy of statistical and machine learning methods for discovery and validation.

Conclusions: While most clinically validated biomarker models derived from omics data have been developed for personalised oncology, first applications for non-cancer diseases show the potential of multivariate omics biomarker design for other complex disorders. Distinctive characteristics of prior success stories, such as early filtering and robust discovery approaches, continuous improvements in assay design and experimental measurement technology, and rigorous multi-cohort validation approaches, enable the derivation of specific recommendations for future studies.

Strengths and limitations of this study

- This scoping review provides an overview of biomarker discovery studies using machine learning analysis of omics data which have led to clinically validated diagnostic and prognostic tools.
- The review discusses shared characteristics of successful biomarker studies as a guidance for study design, discovery and validation method choices for future projects.
- Data extraction and analysis methods focus on deriving recommendations to optimize the design of prospective studies and improve analysis workflows for retrospective studies.
- The review applied minimum eligibility criteria for sample size and statistical validation, but did not assess the quality of the included studies.

Introduction

Personalised medicine is a rapidly developing area in health care research and practice, which aims at providing more effective and safer therapies tailored to the individual patient, by exploiting subject-specific molecular, clinical and environmental data sources (Box 1).

A central tool in personalised medicine and the focus of this study is the machine learning (ML) analysis of omics profiling data to derive molecular biomarker signatures for disease- or drug-based patient stratification (1). The major goals for ML-based omics biomarker development are to develop more reliable and robust tests for drug response prediction, early diagnosis, differential diagnosis or prognosis of the future clinical disease course (2). Omics-derived biomarker signatures may help to guide treatment decisions, and to focus therapies on the right populations to prevent overtreatment, increase success rates, and reduce costs (3). As a research and information tool, they may enable a better monitoring of disease progression and treatment success, and guide new drug development and discovery (4). In contrast to classical single-molecule biomarker approaches, omics signatures have the potential to provide more sensitive, specific and robust predictions of disease-associated outcomes (5).

However, while biomarker discovery projects using omics data have already led to the successful development of clinically validated diagnostic and prognostic tests (6–15), many biomarker studies are discontinued after early development stages or fail in later clinical validation stages. Dedicated statistical and ML methodologies for omics biomarker discovery and validation have been published, as well as recommendations for study design, implementation and reporting (16,17). The distinctive features and approaches which characterize prior successes in translating omics research findings into clinically validated tests have however not yet been investigated in detail. In order to guide future projects on suitable method choices, there is a need for dedicated studies on the key determinants of previous translational successes in ML-based omics biomarker development.

As part of an EU project on “Personalised Medicine Trials” (PERMIT (18)), funded within the H2020 framework, we have therefore investigated the current methodological practices for personalised medicine, covering ML approaches for omics-based patient stratification as a major focus area. While a broader series of questions was established and examined for the overall scoping review (19), for this manuscript, we focused our analysis on biomarker discovery studies that have led to successful, clinically validated FDA-cleared tests or laboratory developed tests (LDTs), to determine their shared and distinctive characteristics compared to studies with no clinical translation. In particular, we aimed to address the following specific research questions:

- Which omics-derived biomarker discovery studies have led to clinically validated tests for patient stratification (LDTs or FDA-cleared tests)?
- What are the key characteristics shared by successful omics biomarker studies and distinguishing them from previously published biomarker studies which have not yet led to clinically validated tests?
- Which types of model building and validation methods have been used to develop clinically validated biomarker signatures, and what are the lessons learned and recommended workflows?
- Which recommendations and guidelines have been proposed to address common challenges in biomarker development using omics data?

1
2
3 These questions lend themselves to a scoping review, because omics-derived biomarker
4 development is still an evolving field, and a preliminary assessment of the potential scope and size
5 of the available biomedical literature on these topics is required as a first step for further follow-up
6 research. Therefore, the objective of this study was to address the above questions by retrieving and
7 examining the current literature on biomarker discovery and validation studies using omics data and
8 ML approaches. While the focus on articles describing discovery and validation approaches covers
9 relevant aspects for clinical translation, we point out that other translational and regulatory aspects,
10 such as the assessment of the clinical efficacy of biomarker-associated treatment decisions, the
11 assessment of cost-effectiveness and research ethics, are not addressed in the present review, but
12 have been discussed in previous dedicated articles (20–24). Our scoping review also does not aim
13 at providing a quantitative benchmark evaluation of different ML approaches, but relevant studies have
14 previously been presented for supervised machine learning (25), and unsupervised clustering (26)
15 and survival prediction (27) on multiple omics data types.

21 **Methods**

22 We conducted a scoping review following the methodological framework suggested by the Joanna
23 Briggs Institute (28). This framework consists of six stages: 1) identifying the research questions, 2)
24 identifying relevant studies, 3) study selection, 4) charting the data, 5) collating, summarising and
25 reporting results, and 6) consultation.

26
27 The scoping review approach was considered most suitable to respond to the broad scope and the
28 evolving nature of the field. Compared to systematic reviews that aim to answer specific questions,
29 scoping reviews present a general overview of the evidence pertaining to a topic and are useful to
30 examine emerging trends, to clarify key concepts and identify gaps (29,30). Before conducting the
31 review, a study protocol was published on the online platform Zenodo (19). Due to the iterative nature
32 of scoping reviews, deviations from the protocol are expected and duly reported when occurred. We
33 used the PRISMA-ScR (Preferred Reporting Items for Systematic reviews and Meta-Analyses
34 extension for Scoping Reviews) checklist to report our results (31) (Online supplementary file 1).

37 *Study identification*

38
39
40 Relevant studies and documents were identified, balancing feasibility with breadth and
41 comprehensiveness of searches. We searched PubMed, EMBASE and Web of Science (last search
42 date: July 27, 2021) for articles describing supervised or unsupervised ML analyses for biomarker
43 discovery or personalised medicine, including both discovery and validation methods. The relevance
44 of the search methodology was ensured by using a strict multi-stage filtering, considering only
45 articles including at least one relevant search term per category from four categories of keywords
46 (“Personalized medicine / Biomarkers”, “Omics”, “Machine Learning” and “Validation”, covering both
47 synonyms for these terms and closely related keywords, see Fig. 1, illustrating the keyword-based
48 search strategy, and Online Supplementary file 2 for the detailed search queries), and subsequently
49 post-filtering the retrieved articles manually to exclude studies not involving omics-based biomarker
50 research or lacking a description of machine learning and validation analyses (see sections on
51 *Eligibility criteria* and *Study selection*). To cover only relevant scientific content, the scope was limited
52 to journal publications and meeting abstracts from international conferences and workshops, and no
53 other grey literature was included. We restricted inclusion to reports published from January 2000 to
54 July 2021 (covering also “online first” articles with official publication date in the future) in English,
55 French, Spanish, Italian and German language. Since to the best of our knowledge, the first clinically
56 validated FDA-cleared omics-derived biomarker signature was published in 2002 (32), only few
57 preliminary discovery studies were expected to have taken place significantly earlier than 2002, and
58 we therefore did not extent the search further backwards in time than January 2000.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Eligibility criteria

We included peer-reviewed methodology articles, review articles, opinion articles on supervised and unsupervised ML methods for omics disease prediction and stratification and associated statistical cross-validation and multi-cohort validation methods (addressing accuracy, robustness, and clinical relevance). Only approaches tested on real-world biomedical omics data were reviewed, while studies relying purely on simulated data or were excluded. We also excluded papers on biomarker methods without a demonstrated biomedical application, and those with insufficient sample size (i.e., removing studies covering less than 50 samples per group for the main conditions studied, unless a dedicated power calculation was presented) or statistical validation (i.e., lack of clear descriptions of cross-validation or external testing methodology, performance metrics and test statistics). These exclusion criteria were not specified in the generic review protocol, but they were agreed among the authors prior to the screening process.

To cover both data from original research papers and prior systematic reviews, we extracted information from three main article types: (1) applied research papers, (2) methodology articles with demonstrated applications, and (3) review articles on methods, applications and validation approaches.

Apart from these inclusion and exclusion criteria, for the final result presentation, the statistical investigations covered all selected articles, whereas the detailed discussion of study characteristics focused on the studies that led to clinically validated biomarker signatures tested on multiple cohorts with large sample sizes (i.e., studies using a power calculation to demonstrate the adequacy of the chosen sample sizes, or covering hundreds or thousands samples per studied subject group).

Study selection

We exported the references retrieved from the searches into the online tool Rayyan (33). Duplicates were removed automatically using the reference manager Endnote X9 (Clarivate Analytics, Philadelphia, United States) and manually by the reviewers. One reviewer loaded the retrieved records into the online screening tool Rayyan (33), and two reviewers confirmed the eligibility independently by covering both the screening for all records and the full-text review for the articles pre-selected by the screening. Disagreements were solved by consensus.

Charting the data and synthesis of results

We designed a data extraction form using Excel (Online supplementary file 3). General study characteristics extracted covered author names, title, citation, type of publication (e.g., journal article, meeting abstract), study population and sample size (if applicable), methodology/study design, and outcome measures (if applicable). Specific items associated with the topic of the scoping review included the study type (e.g., Case-control study, differential diagnosis study, prognostic study, review – methods, review – applications, review – validation); the article type (journal or conference article), the generic ML domain (e.g., supervised/unsupervised); and the name of specific approaches for outcome prediction and for validation. Moreover, to capture key findings related to the review questions, relevant sentences were extracted from each reviewed article, and if needed, complemented by a brief explanatory remark, and by writing out abbreviations used in the original text.

1
2
3 The reviewers piloted the data extraction form using five records from the retrieved article collection.
4 Two reviewers (EG, AR) working independently extracted the data from the included articles. In the
5 case of disagreements, consensus was obtained by discussion.
6

7 In the final full-text review stage, the pre-selected articles were grouped by topic, categorizing articles
8 into applied vs. methodological studies, supervised vs. unsupervised analyses, and assigning
9 algorithm type identifiers to each article (review articles and papers on validation methodologies
10 were considered as separate categories without a specific algorithm type assignment). The full-text
11 review and categorization of articles into different publication types was done through independent
12 manual inspection by the two reviewers.
13

14 While the information on sample sizes and validation methods was documented as part of the data
15 extraction (Online supplementary file 4, a spreadsheet version has been made available on the
16 online platform Zenodo (34)), it was not within the remit of this scoping review to assess the
17 methodological quality of individual studies included in the analysis.
18
19

20 21 22 *Consultation exercise*

23 The members of the PERMIT consortium, associated partners, and the PERMIT project Scientific
24 Advisory Board discussed the preliminary findings of the scoping review in a 2-hour online workshop.
25
26

27 28 *Patient and public involvement*

29 The European Patients' Forum is a member of PERMIT project. Although not directly involved in the
30 conduction of the scoping review, they received the draft review protocol for collecting comments
31 and feedback.
32
33

34 35 36 37 **Results**

38 39 *Study selection and general characteristics of reports*

40 We retrieved 1563 abstracts from the literature search. After the removal of duplicates, we screened
41 the remaining 1475 abstracts for eligibility. 619 records were excluded, while 856 abstracts were
42 retained for the full-text assessment. Finally, we included 352 articles that passed all filtering criteria
43 in the data extraction and analysis (see flow chart in Fig. 2, and Online supplementary file 4, providing
44 the reference for each selected article, as well as information on the study type and methodology,
45 the outcome measures, the validation type, and representative sentences from each article on the
46 main study results and key findings; a spreadsheet version of this table has been available on the
47 online platform Zenodo (34)).
48
49

50 The full-text article review revealed that many studies did not meet the pre-defined inclusion criteria:
51 371 articles (43%) were removed because of an insufficient sample size, and 105 further articles
52 (12%) were excluded because they provided insufficient details on the validation results or
53 methodology (see Fig. 2). This shows that the challenges of recruiting an adequate number of
54 participants per study group or conducting sufficient omics profiling experiments for robust model
55 building and validation are not met in a large proportion of omics biomarker studies. Moreover, many
56 studies lack adequate documentation for the study design and validation.
57
58

59 For the selected articles that cover primary research on omics biomarker studies, the majority (78%)
60 rely entirely on an internal validation involving data from only a single cohort, whereas studies that

1
2
3 use an external validation on an independent cohort are still underrepresented (only 12% of articles
4 describe both an internal cross-validation and an external cohort validation, and an additional 10%
5 include an external validation, but do not report internal cross-validation results). However, when
6 comparing the numbers of published studies over different periods of time during the past 20 years,
7 the relative proportion of studies including an external validation has increased in recent years (see
8 Fig. 3), suggesting a growing recognition of the importance of independent, multi-cohort validation.
9

10 Next, we investigated the countries of origin for the selected articles, showing that the United States
11 of America (USA) are contributing the largest proportion of validated biomarker studies (28%),
12 followed by China (18%), Canada (5%), Germany (4%), and the United Kingdom and India (both
13 3%; see also Fig. 4, providing a map visualization of the country statistics). These country
14 representations show limited correlation with population sizes and may largely reflect worldwide
15 variation in relative biomedical research productivity reviewed in previous study (35). Since the most
16 prolific countries in the development of molecular diagnostics have already set up policies and
17 regulations for omics- and ML-based *in vitro* diagnostics and medical devices (e.g., see the life cycle
18 regulation of AI- and ML-based software devices in the USA (36)), they may also provide a role
19 model for countries still in the process of establishing similar regulatory frameworks.
20
21

22 When inspecting the representation of study design types in the filtered article collection, the great
23 majority of documents described diagnostic studies (67%), prognostic and survival prediction studies
24 were covered in 8% of articles, and studies examining therapy or drug response in 7% (see Fig. 5).
25 Apart from this, 13% of articles were reviews on methodologies and applications in the field, and 5%
26 of articles described other rare study types (e.g. tissue-of-origin prediction studies or combinations
27 of different study types).
28
29

30 Since a detailed discussion of all filtered articles is not within the scope of the present review, in the
31 following, we focus on reviewing representative omics biomarker studies which achieved the highest
32 validation level, i.e., clinical approval of the developed molecular signature as an LDT or FDA
33 approved test (see the overview of studies in Table 1 and the FDA web-site (37)). We investigate
34 the shared features of these successful studies, examine how they address common shortcomings
35 and missing features of other reviewed studies, and summarize the lessons learned.
36
37
38
39
40

41 *Success stories in omics-based biomarker signature development*

42 *Cancer approved omics-derived diagnostic tests (9 studies)*

43
44
45 The first and most well-known omics-derived molecular test to receive FDA clearance was
46 *MammaPrint*[®], a prognostic signature using the RNA expression activity of 70 genes to estimate the
47 risk for distant tumours metastasis and recurrence in early-stage breast cancer patients (6,32,38–
48 41). This test was developed at the Netherlands Cancer Institute, using DNA microarray analysis to
49 investigate primary breast tumours of 117 patients. Supervised ML was applied to the resulting data
50 to identify a highly predictive gene signature for a short interval to distant metastases in lymph node
51 negative patients (32).
52
53
54

55 A distinctive feature of the development approach behind this signature in comparison to other
56 reviewed studies was the multi-stage filtering and cross-validation strategy used in the initial
57 discovery study, which may explain the repeated confirmation of the signature in later validation
58 studies (6,38–41). From 25k genes represented on the DNA microarrays, only those significantly
59 regulated in more than 3 tumours out 78 sporadic lymph-node negative patients were preselected,
60 and further filtered by retaining only the genes with a minimum absolute correlation with the disease

1
2
3 outcome of 0.3. The resulting list of 231 genes, rank-ordered by absolute correlation, was
4 investigated by sequentially adding the next top 5 genes from the list to a candidate ML classifier
5 and evaluating its performance by leave-one-out cross-validation (LOOCV). This procedure was
6 repeated as long as the estimated accuracy of the classifier improved, providing a final candidate
7 signature of 70 genes. The final signature was validated on multiple independent test sets, including
8 a set of 19 external samples in the original study and several additional validations on independent
9 cohorts in follow-up studies (6,38–41).

11 The *MammaPrint* signature provided the role model for the subsequent development of a similar
12 prognostic test for colon cancer, *ColoPrint*[®] (42–47). This test aims at detecting the approx. 20% of
13 patients with stage II colon cancer expected to experience a relapse and develop distant metastases.
14 It uses an 18-gene expression signature, developed by analysing DNA microarray data in a similar
15 manner to the *MammaPrint* approach. The diagnostic approach has been commercialized as an
16 LDT to assist physicians in selecting treatment options for colon cancer patients. Similar to
17 *MammaPrint*, the signature development was characterized by extensive discovery and validation
18 studies, which involved multiple statistical reproducibility, stability and precision analyses for
19 independent, large-scale patient cohorts (48).

22 Another widely used cancer-related LDT, which received FDA clearance in 2013, is the *Prosigna*[®]
23 *Breast Cancer Prognostic Gene Signature Assay*, previously called *PAM50* test (49–53). This assay
24 assesses mRNA expression for a signature of 58 genes (50 target genes + 8 endogenous control
25 genes) to predict the risk of distant recurrence for hormone-receptor-positive breast cancer between
26 5 to 10 years after diagnosis (prerequisites are that the patients have been treated with hormonal
27 therapy and surgery, and are stage I or stage II lymph-node negative, or in stage II with one to three
28 positive nodes). The test development started with a microarray discovery study and involved a
29 multistage filtering, using consecutive applications of statistical tests and cross-validation to propose
30 a subset of candidate gene markers (54). The authors compared the reproducibility of classification
31 scores obtained with these markers for three centroid-based prediction methods to ensure the
32 robustness of the methodology. By further developing the approach into a more sensitive PCR-based
33 test, and later into an assay using the NanoString nCounter Dx Analysis System, the predictive
34 performance was improved in a step-wise fashion. The original discovery study was characterized
35 by significantly larger sample sizes than the majority of reviewed biomarker studies, with a training
36 set of 189 samples, test sets of 761 patients evaluated for prognosis, and 133 patients evaluated for
37 prediction of pathologic complete response to treatment with taxane and anthracycline. These study
38 design features in combination with multi-stage filtering and validation approaches, and improved
39 measurement technology during the course of the study, may explain the successful progression of
40 the *PAM50* test to FDA clearance. The test has only three genes in common with the *MammaPrint*
41 approach (*KNTC2*, *MELK*, *ORC6L*), which may be explained by the different technical and analytical
42 approaches used, but a previous comparative evaluation concluded that the tests provide broadly
43 equivalent risk information for females with oestrogen receptor (ER)-positive breast cancers (55).

48 Among the LDTs for breast cancer prognosis, *Oncotype DX*[®] is a further test commonly used in
49 clinical practice (8,56–59). The underlying gene signature consists of 16 cancer-associated genes
50 and 5 reference genes, and is therefore often also referred to as ‘21-gene assay’. Its main application
51 is to predict risk of recurrence in oestrogen-receptor positive tumours. The relevance of this
52 prognostic tool for treatment selection may be explained by the strong association of the provided
53 recurrence score with the probability of positive treatment response to chemotherapy (60). *Oncotype*
54 *DX* was developed using a consecutive refinement procedure, starting with the RT-PCR assessment
55 of 250 candidate genes across 447 patients from three distinct studies to identify the 21-gene
56 signature after multiple filtering steps. A recurrence score algorithm built using the signature as input
57 was clinically validated on 668 independent patients (61). The selection of the 16 cancer-related
58 genes included in the assay involved scoring the performance of the candidate features in all three
59 studies and the consistency of the primer/probe performance in the assay (62). Thus, particular
60

1
2
3 strengths of the development process for this LDT include the consideration of both technical
4 robustness and statistical robustness of the assay across distinct cohorts. The *Oncotype DX*
5 signature shares one gene with *MammaPrint (SCUBE2)*, and 9 genes with the Prosigna / PAM50
6 test (*BIRC5, CCNB1, MYBL2, MMP11, GRB7, ESR1, PGR, BCL, BAG1*). However, an independent
7 clinical validation of Oncotype DX and the PAM50 signature for estimating the likelihood of distant
8 recurrence in ER-positive, node-negative, post-menopausal breast cancer patients treated with
9 endocrine therapy suggested that the PAM50 signature provided more prognostic information than
10 Oncotype DX (63).
11

12
13 While the first validated omics biomarker signatures were developed for breast cancer, similar
14 diagnostic and prognostic tools have followed for other cancer types. One of these is the Decipher®
15 Prostate Cancer Test (9,64–68), which differs from other omics-derived diagnostic tools by being
16 provided together with a software platform and database, the Decipher Genomic Resource
17 Information Database (GRID), that captures 1.4 million expression markers per patient to facilitate
18 personalised care. The test itself uses 22 preselected RNAs to predict clinical metastasis and
19 cancer-specific mortality for patients who have undergone radical prostatectomy. An initial discovery
20 study by the Mayo Clinic (Rochester, MN, USA) investigated a cohort of 545 such patients, split into
21 a training (n = 359) and a validation cohort (n = 186). Similar to other LDTs, the discovery started
22 with a genome-wide profiling and used both statistical and ML analyses for filtering. First, *t*-tests
23 were applied (reduction from 1.4 mil. to 18,902 differentially expressed RNAs), then regularized
24 logistic regression (reduction to 43 candidate markers), and finally a random forest-based feature
25 selection (reduction to final set of 22 RNAs). Apart from testing the signature in the validation cohort,
26 further external validations were performed in subsequent studies (9,64–68). Overall, distinctive
27 strengths of the used approach include the improved interpretability of the test results through
28 supporting analyses on the GRID platform, and the robustness of the discovery and validation
29 approach, involving large sample sizes and several complementary statistical and ML assessments.
30
31

32
33 While most diagnostic tests in oncology have been designed for specific cancer types, a dedicated
34 LDT has also been developed for cancers of unknown or uncertain diagnosis. The Cancer Type ID®
35 test by bioTheranostics distinguishes between 50 different tumour types using a 92-gene RT-PCR
36 expression measurement signature (15,69–71). This signature was derived from analyses of a
37 microarray data collection covering 446 frozen tumour samples and 112 formalin-fixed, paraffin-
38 embedded (FFPE) samples of both primary and metastatic tumours. Modelling steps involved *k*-
39 nearest neighbour clustering and classification, and a genetic algorithm to explore the search space
40 of possible feature subset selections. After successful cross-validation (84% accuracy) and external
41 validation (82% accuracy on 112 independent FFPE samples), the microarray-based signature was
42 further developed to use more sensitive RT-PCR measurements. Testing the new approach on an
43 independent validation set provided an increased accuracy (87%). Distinctive characteristics of the
44 development process that may have contributed to the positive validation include the efficient and
45 extensive exploration of the search space of possible gene subset selections via a genetic algorithm,
46 the large sample sizes used for discovery and validation, and the transfer of the assay from
47 microarrays to the more sensitive RT-PCR platform.
48
49

50
51 The first omics-derived biomarker signatures addressed only the most frequent cancer types, but
52 more recent applications in oncology focus on the diagnosis of less common malignancies, such as
53 thyroid cancer. Typically, deciding whether a thyroid nodule is benign or cancerous is possible via a
54 fine needle aspiration (FNA) biopsy, without requiring more complex measurements or analyses.
55 However, while direct FNA-based diagnosis is feasible in most cases, indeterminate results can
56 occur (72). To help prevent unnecessary surgeries for the corresponding patients, a molecular
57 signature and LDT known as the Afirma™ Gene Expression Classifier (GEC) has been developed
58 to discriminate benign from cancerous thyroid nodules (72–77). The original discovery study behind
59 the GEC signature used mRNA expression analysis in 315 thyroid nodules, covering 178
60 retrospective surgical tissues and 137 prospectively collected FNA samples. Two ML classifiers were

1
2
3 trained separately on surgical tissues and FNAs, assessing the test set performance on 48
4 independent, prospective FNA samples (50% of which had indeterminate cytopathology).
5 Discriminative features were selected using a linear modelling approach implemented in the software
6 Limma, and a linear support vector machine was applied for model building and performance
7 estimation via 30-fold cross-validation (CV). The successful cross-validation results were confirmed
8 on multiple distinct cohorts (72,75–78). While the internal validation used in the initial study cannot
9 address cohort-specific biases, the combined use of established feature selection and modelling
10 approaches, and the subsequent external validation across multiple cohorts with large sample sizes
11 may account for the successful translation of this signature.
12
13

14 Most omics-based diagnostic tests identified in our study rely purely on gene expression profiling
15 data. However, more recently, first multi-omics signatures for diagnostic purposes have been
16 developed. One of the first LDTs that integrated information from both RNA and DNA sequencing
17 was the FoundationOne® Heme assay (14,79–81). This assay aims to detect hematologic
18 malignancies, sarcomas, pediatric malignancies, or solid tumours (including among others
19 leukaemias, myelodysplastic syndromes, myeloproliferative neoplasms, lymphomas, multiple
20 myeloma, Ewing sarcoma, Leiomyosarcoma, and paediatric tumours). The test identifies four types
21 of genomic alterations (base substitutions, insertions and deletions, copy number alterations,
22 rearrangements) and reports microsatellite instability and tumour mutational burden to facilitate
23 clinical decision making. This approach was originally developed and evaluated using reference
24 samples of pooled cell lines in order to model the main characteristics that determine the test
25 accuracy, including mutant allele frequency, indel length and amplitude of copy change (79). A first
26 validation using 249 independent FFPE cancer samples, which had already been characterized by
27 established assays, confirmed the accuracy of the test. External validation studies on independent
28 cohorts corroborated the utility of the test for further diagnostic applications (14,82). The study results
29 highlight the potential of integrating diverse biological data sources in order to obtain more robust
30 and reliable predictions, a strategy that may be promising in particular for complex disorders that
31 involve very heterogeneous phenotypes.
32
33
34

35 A common limitation of genomic profiling approaches for diagnostic testing is that most analyses
36 have to be performed in centralized specialty laboratories, which limits a wider use and results in
37 long waiting times. To address this shortcoming, the Elio™ Tissue Complete assay, an *in vitro*
38 diagnostic test cleared in 2020 by the FDA for assessing somatic mutations and tumour mutation
39 burden (TMB) in solid tumours, has been developed as an integrated DNA-to-report approach to
40 enable a decentralized evaluation in all diagnostic labs with next generation sequencing (NGS)
41 technology (83). The analytical performance of the test was assessed by comparing it with the
42 FoundationOne test (see above) using a concordance analysis on 147 tumour specimens. It
43 provided a positive percent agreement (PPA) above 95% for single nucleotide variants (SNVs) and
44 insertions/deletions, and 80–83% PPA for copy number alterations and gene translocations (83). The
45 test has recently also been applied to investigate the response to immune checkpoint inhibitors (ICI)
46 in metastatic renal cell carcinoma (mRCC), using a retrospective evaluation of SNVs, TMB,
47 microsatellite status and genomic status of antigen presentation genes (84). While no correlation
48 between treatment response and TMB was observed, one third of patients with progressive disease
49 following ICI therapy displayed loss of heterozygosity of major histocompatibility complex class I
50 genes (LOH-MHC) vs. 6% of disease control patients, suggesting that loss of antigen presentation
51 may restrict ICI response (84). In summary, the Elio Tissue Complete assay provides an example of
52 how integrating NGS analyses with bioinformatics in a combined DNA-to-report approach could help
53 to broaden the access to genomic diagnostics for both clinical and research applications.
54
55
56
57
58
59

60 *Non-cancer approved omics-derived diagnostic tests (4 Studies)*

1
2
3 While most clinically approved omics-derived diagnostic tests have been developed in the field of
4 oncology, one of the first LDTs that received FDA clearance for a non-cancer disease was the
5 AlloMap® Heart test (13,85–87). It uses a gene expression signature of 11 target genes and 9 control
6 genes in peripheral blood from heart transplant recipients to estimate the risk for acute cellular
7 cardiac allograft rejection. The development process involved statistical analyses of leukocyte
8 microarray profiling data from 285 samples, and subsequent RT-PCR validation and bioinformatics
9 post-processing (13). Prior knowledge from database and literature mining was included in the
10 analysis by mapping the data to known alloimmune pathways. This allowed the researchers to
11 narrow down 252 candidate marker genes. An RT-PCR validation on 145 samples confirmed 68 of
12 these candidate genes, which distinguished rejection samples from quiescent samples according to
13 a T-test ($p < 0.01$). Six genes were eliminated due to significant variation in gene expression with
14 sample processing time. Next, the investigators averaged correlated gene expression levels to
15 create robust meta-level features, called ‘metagenes’, and added 20 of these features as new
16 variables. A linear discriminant analysis was applied, providing a prediction model using four
17 individual genes and three metagenes, which aggregate information from 11 original genes. Finally,
18 bootstrap validation procedures and external test set validations were performed to confirm the
19 accuracy of this signature. Overall, distinctive aspects of the development approach for the AlloMap
20 signature include the knowledge-based gene discovery, a comprehensive RT-PCR validation of
21 candidate genes, and the robust bootstrap and external validation analyses.

22
23
24
25 The first clinically validated LDT for a cardiovascular indication derived from omics data was the
26 Corus® CAD test, developed to identify coronary artery disease (CAD) in stable non-diabetic patients
27 (11,88–91). In contrast to most other omics-based tests, Corus CAD is not a pure molecular
28 signature test, but takes the clinical covariates gender and age into account. The initial discovery
29 study used a retrospective microarray analysis of blood samples from 195 diabetic and non-diabetic
30 patients from the Duke University CATHGEN registry. After ranking the studied genes by the
31 statistical significance of group differences and prior biological knowledge on their disease
32 relevance, 88 genes were selected for RT-PCR validation. Because diabetes status as a clinical
33 covariate was significantly associated with the observed gene expression alterations, and the
34 identified CAD-associated genes did not overlap between diabetic and non-diabetic patients, the
35 authors decided to limit follow-up work to non-diabetic patients. In a prospective clinical trial,
36 microarray profiling was conducted on blood samples from 198 patients, and top-ranked genes were
37 further validated using RT-PCR for 640 blood samples. After multiple filtering steps, taking into
38 account statistical significance in T-tests, biological relevance, gene correlation clustering and cell-
39 type analyses, a final signature of 23 genes was derived, composed of 20 CAD-associated genes
40 and 3 reference genes (92). To maximize the predictive performance, the final prediction algorithm
41 was optimized to adjust for differences associated with age and gender. Compared to most other
42 reviewed studies, the Corus CAD approach stands out by taking clinical covariates into account in
43 the final prediction model, including an intermediate critical review and adjustment of the inclusion
44 criteria (limiting the focus to nondiabetic patients), and integrating complementary filtering and
45 validation analyses on large sample sizes.

46
47
48
49 For inflammatory diseases, a first omics-derived signature recently received approval for measuring
50 rheumatoid arthritis (RA) inflammatory disease activity, the Vectra® DA multi-biomarker test (93–97).
51 It uses blood serum samples and multi-spot 96-well immunoassay plates to assess serum
52 concentrations of 12 protein biomarkers associated with the pathobiology of RA. The original Vectra
53 DA score, which combines these measurements into a composite score between 1 and 100, was
54 assessed via multivariate regression and displayed a high predictive power in estimating a standard
55 RA score, the Disease Activity Score in 28 joints using the C-reactive protein level (DAS28-CRP), in
56 both seropositive (AUC 0.77, $P < 0.001$) and seronegative (AUC 0.70, $P < 0.001$) patients (97). This
57 score was later adjusted for age, gender and adiposity (based on leptin concentration), and validated
58 in two cohorts against DAS28-CRP as a prognostic test for radiographic progression during the next
59 year. The results showed that the new adjusted score was the most accurate, independent predictor
60

of progression, with the rate of progression increasing from < 2% in the low (1-29) adjusted score category to 16% in the high (45-100) category (95). Overall, the Vectra DA approach illustrates the utility of omics-based biomarker signatures for prognostic applications in inflammatory disorders, and further highlights the benefit of integrating omics signatures with information from clinical covariates.

For neurodegenerative disorders, clinically approved diagnostic and prognostic omics-derived tests are still lacking. However, recently the Helix® Genetic Health Risk App for Late-onset Alzheimer's Disease (AD) was cleared by the FDA for over-the-counter use. It detects clinically relevant variants in genomic DNA isolated from human saliva of individuals ≥ 18 years in order to report and interpret genetic health risks, and evaluates the information of variants with established genome-wide significant associations to AD. When tested on 99 human saliva samples, the accuracy was 100% with a lower 95% CI bound of 96.3% (98). The approach uses a whole exome sequencing (WES) constituent device, the Helix® Laboratory Platform (99–101), as a qualitative *in vitro* diagnostics approach covering measurements for approximately 20k genes. The Helix Laboratory Platform has received FDA clearance through a new regulatory approval pathway established by the FDA for WES devices (Regulation 21 CFR 866.6000). Due to the generic applicability of the WES profiling assay used by this platform, called Exome+, the assay has also been applied to find statistically significant gene-based associations for several other phenotypes in large-scale cohort studies (99) and to identify carriers of autosomal dominant diseases by population-based genetic screening (101). Thus, the Helix Laboratory Platform provides a first example for a new approval pathway for omics-based diagnostic tests, in which a clinically approved genomic testing device is not anymore linked to a single diagnostic application or a specific disease type. Instead, the market authorization for diagnostic tests is obtained separately from the device and facilitated and accelerated by the prior approval of the constituent measurement device. For the future development of omics-derived biomarker signatures, this may allow researcher to focus on demonstrating the clinical utility of a new signature, while the analytical validity of the underlying testing device has already been established previously.

Discussion

Statement of principal findings

The scoping review of articles on patient stratification using omics data revealed common limitations in the study design for many published biomarker development projects, such as insufficient and imbalanced sample sizes per study group and inadequate validation methods, but also identified multiple studies that have led to validated diagnostic and prognostic tests. These success stories were investigated in more detail to identify common characteristics in the study design, discovery and validation methods, which may have supported the clinical translation of the initial findings. Fig. 6 outlines key shared aspects that are possible determinants of the study success and could help to guide future biomarker investigations. In particular, they cover the following main features:

(1) A sample size selection, study group and replicate design that provides adequate statistical power for the ML analyses;

2) The application of robust statistical filtering and evaluation schemes (including multiple layers of statistical and ML-based feature selection, combined statistical and biological filters, robust validation schemes that involve multiple cross-validation, bootstrapping and external validation analyses, using multiple suitable and complementary performance metrics, and providing information on the statistical variation and confidence intervals for the performance estimates, see Fig. 7 for an overview of recommended generic steps for robust model building and evaluation);

3) Clarity of the study scope and goals (involving clear inclusion and exclusion criteria, primary and secondary outcomes, and decision processes to make necessary adjustments due to new knowledge gained during the project, such as the adjusted inclusion criteria in the Corus CAD study and the progression from non-targeted microarray technology to higher-sensitivity RT-PCR in the case of the Prosigna test and the Cancer Type ID test);

4) Completeness and reproducibility of the study documentation (covering details on used instruments, parameters and settings, reproducible methods descriptions, and information on data provenance);

5) Interpretability and biological plausibility of the created predictive models (including explainable and justifiable predictions, human-interpretable model descriptions, and biologically plausible models that agree with the current mechanistic understanding of the studied disorder);

6) Integration of prior biological knowledge into the predictive feature selection, model building and validation procedures (e.g., using public data on disease-associated molecular pathways and networks; complementary clinical and real-world data, and relevant multi-omics data).

Strengths and Limitations

The majority of methodological recommendations derived from the study relate to the early planning and study design for biomarker discovery projects, involving considerations associated with the choice of the study group, sampling and blocking design, the measurement technology, and the input and output variables (16,17). These recommendations are therefore mainly applicable to prospective studies. For retrospective biomarker investigations of already collected data, the suggestions derived from the review are limited to guidance on improving analysis workflows, e.g. for filtering and evaluation analyses, the integration of prior knowledge from multi-omics data and public annotation databases, and the choice of robust and interpretable modelling approaches for the generation of biologically plausible and reproducible prediction models. While the focus of the review on studies that have already led to validated biomarker models and that fulfil minimum requirements for sample size and statistical model assessment helps to ensure the quality of the selected articles, no further quality evaluation was performed. The reader should also note the generic limitations of machine learning methods which can affect all biomarker studies: These include the necessity for a representative coverage of the relevant outcomes in the training and validation groups, a sufficiently comprehensive and sensitive coverage of informative predictor variables in the data for the outcomes of interest, which may not be achievable for omics data from tissues and body fluids with limited disease relevance or measurement sensitivity, and a sufficient data quality in terms of the influence of systematic biases and noise. Moreover, for multi-omics biomarker analyses, in addition to adequate pre-processing and machine learning approaches, suitable strategies and methods for the integration of diverse omics data are also needed. These multi-omics data integration strategies were not within the scope of the present review, but have been reviewed in previous publications (102–104). Finally, more recent methodological developments in the machine learning and cross-validation analysis of omics data, such as meta-learning (105) and bolstered cross-validation (106), have only limited coverage among the articles that passed the eligibility criteria, and will therefore require further dedicated study in the future.

Discussing important differences in results

Previous reviews of ML approaches using omics data for patient stratification have focused on domain-specific analyses for specific types of diseases, or specific types of ML methodologies (107–115). By contrast, this scoping review focuses on disease-agnostic workflows with generic applicability across complex human disorders involving multifactorial molecular alterations. The

1
2
3 coverage of statistical and ML approaches for stratification does not aim to provide a detailed
4 discussion of specific algorithms, statistical methods or scoring metrics, but rather at identifying key
5 determinants of success for generic analysis and validation workflows in biomedical stratification
6 studies. Therefore, the results describe general workflow characteristics that distinguish omics
7 biomarker studies with clinical translation from other studies, and cover associated disease-agnostic
8 recommendations for future studies, whereas method recommendations specific to particular
9 disease types or ML analysis types are covered elsewhere in domain-specific reviews (107–115).
10
11
12

13 *Meaning of the study: implications for clinicians and policymakers*

14
15 The previous clinical translation successes in omics-based biomarker development reviewed in this
16 study, which have mostly been achieved in the field of oncology, highlight the potential for developing
17 similar biomarker signatures for further disease indications. In contrast to conventional statistical
18 biomarker discovery approaches, which focus on identifying single-molecule markers, systems-level
19 analysis of omics data using multivariate ML approaches can identify multifactorial signatures which
20 are robust against noise in individual gene or protein measurements, and more biologically insightful
21 by reflecting disease-associated cellular process alterations in a more comprehensive fashion.
22
23

24 This scoping review has identified common characteristics of omics studies which have led to
25 clinically validated diagnostic and prognostic tests. Thus, the conclusions drawn on recommended
26 practices for sample size selection, biological data filtering and ML, and the implementation of
27 adequate validation schemes may help to guide clinical researchers on study design choices and
28 the selection of analysis methodologies. Additionally, the scoping review results can help to raise
29 awareness of common pitfalls, such as issues associated with batch effects, biases, confounding
30 factors, lack of statistical power, and multiple hypothesis testing, and thus contribute to preventing
31 these failure causes in biomarker development. For policymakers and funding bodies, findings on
32 the distinctive characteristics of studies with successful clinical biomarker translation, e.g.
33 concerning the specific requirements for robust cross-validation and external result validation
34 methods, may provide relevant information for the design of public and private funding schemes for
35 biomedical research. Risks in funded research projects may be addressed upfront through
36 appropriate guidelines and regulations for the study design and validation (e.g. recommendations
37 on power calculations and specific validation and documentation requirements). Finally, the scoping
38 review results can guide clinicians involved in biomarker discovery on how to make better use of
39 available public knowledge and data sources, e.g. cellular pathway and molecular interaction
40 databases, that may allow them to exploit prior knowledge effectively, and create more robust and
41 interpretable biomarker models.
42
43
44
45
46

47 *Unanswered questions & future research*

48 Since the recommendations and guidelines identified from the reviewed articles are mostly derived
49 from established biomarker discovery and validation approaches, new methodologies and upcoming
50 trends could only be covered to a limited extent and may lead to changed recommendations in the
51 future. In particular, in the reviewed patient stratification studies, some of more recently introduced
52 ML concepts (e.g. transfer learning, distance metric learning, semi-supervised learning, structured
53 machine learning, meta learning, multi-view learning, and generative models), data processing
54 techniques (e.g. new dimension reduction approaches, outlier removal methods, data augmentation
55 techniques), and model validation methods (e.g. bootstrapping or bolstered cross-validation,
56 uncertainty quantification), are still underrepresented among the eligible studies reviewed, and may
57 provide suitable topics for follow-up research.
58
59
60

1
2
3 Overall, while the currently available literature on validated stratification biomarkers already provides
4 ample information on common pitfalls and established practices, the development of widely accepted
5 standard guidelines on methodologies for omics biomarker discovery will require further knowledge
6 exchange and deliberation among stakeholders in the field. In particular, integration of domain-
7 specific expertise in discussions involving clinicians, experimental and data scientists, and regulatory
8 and legal experts is required as a follow-up effort to derive comprehensive methodological guidelines
9 for future biomarker development.
10
11
12
13
14

15 **Acknowledgments**

16
17 The authors thank Vanna Pistotti for her assistance with search strategy development and
18 conduction.
19
20
21

22 **List of Figures**

23 Figure 1: Keyword based search strategy for the scoping review

24 Figure 2: Study selection flow diagram

25 Figure 3: Validation methods used in omics biomarker studies

26 Figure 4: Map representation of country statistics for the selected articles

27 Figure 5: Representation of study types among the selected articles

28 Figure 6: Characteristics of successful omics-based studies

29 Figure 7: Recommended generic workflow for biomarker development using machine learning
30 analysis of omics data
31
32
33
34
35
36
37
38
39
40

41 **Definitions (In boxes)**

42 **Box 1: *What is Personalised Medicine?***

43
44
45 According to the European Council Conclusion on personalised medicine for patients, personalised
46 medicine is 'a medical model using characterisation of individuals' phenotypes and genotypes (e.g.,
47 molecular profiling, medical imaging, lifestyle data) for tailoring the right therapeutic strategy for the right
48 person at the right time, and/or to determine the predisposition to disease and/or to deliver timely and
49 targeted prevention (116).

50
51 In the context of the Permit project, we applied the following common operational definition of
52 personalised medicine research: a set of comprehensive methods (methodology, statistics, validation,
53 technology) to be applied in the different phases of the development of a personalised approach to
54 treatment, diagnosis, prognosis, or risk prediction. Ideally, robust and reproducible methods should cover
55 all the steps between the generation of the hypothesis (e.g., a given stratum of patients could better
56 respond to a treatment), its validation and pre-clinical development, and up to the definition of its value in
57 a clinical setting (19).
58
59
60

References

1. Reel PS, Reel S, Pearson E, Trucco E, Jefferson E. Using machine learning approaches for multi-omics data analysis: A review. *Biotechnol Adv.* 2021;49(107739).
2. Goecks J, Jalili V, Heiser LM, Gray JW. How Machine Learning Will Transform Biomedicine. Vol. 181, *Cell.* 2020. p. 92–101.
3. Jiang Y, Wang M. Personalized medicine in oncology: Tailoring the right drug to the right patient. *Biomarkers in Medicine.* 2010.
4. Hopp WJ, Li J, Wang G. Big Data and the Precision Medicine Revolution. *Prod Oper Manag.* 2018;27(9):1647–64.
5. Glaab E. Using prior knowledge from cellular pathways and molecular networks for diagnostic specimen classification. *Brief Bioinform.* 2016;
6. Cardoso F, van't Veer LJ, Bogaerts J, Slaets L, Viale G, Delaloge S, et al. 70-Gene Signature as an Aid to Treatment Decisions in Early-Stage Breast Cancer. *N Engl J Med.* 2016;375(8):717–29.
7. Bachleitner-Hofmann T, Simon I, Salazar R, Tabernero J, Rosenberg R, van der Akker J, et al. Development and Validation of a Robust Molecular Diagnostic Test (COLOPRINT) for Predicting Outcome in Stage II Colon Cancer Patients. *Ann Oncol.* 2012;
8. Knezevic D, Goddard AD, Natraj N, Cherbavaz DB, Clark-Langone KM, Snable J, et al. Analytical validation of the Oncotype DX prostate cancer assay - a clinical RT-PCR assay optimized for prostate needle biopsies. *BMC Genomics.* 2013;14(1).
9. Torres A, Alshalalfa M, Tomlins SA, Erho N, Gibb EA, Chelliserry J, et al. Comprehensive Determination of Prostate Tumor ETS Gene Status in Clinical Samples Using the CLIA Decipher Assay. *J Mol Diagnostics.* 2017;19(3):475–84.
10. Angell TE, Babiarz J, Barth N, Blevins T, Duh Q, Ghossein RA, et al. Clinical validation of the AFIRMA genomic sequencing braf V600E classifier. *Thyroid [Internet].* 2017;27:A50. Available from: <http://www.embase.com/search/results?subaction=viewrecord&from=export&id=L624116485>
11. Ladapo JA, Budoff MJ, Sharp D, Kuo JZ, Huang L, Maniet B, et al. Utility of a Precision Medicine Test in Elderly Adults with Symptoms Suggestive of Coronary Artery Disease. *J Am Geriatr Soc.* 2018;66(2):309–15.
12. Tabari E, Lovejoy AF, Lin H, Bolen CR, Saelee SL, Lefkowitz JP, et al. Molecular characteristics and disease burden metrics determined by next-generation sequencing on circulating tumor DNA correlate with progression free survival in previously untreated diffuse large B-cell lymphoma. *Blood [Internet].* 2019;134. Available from: <http://dx.doi.org/10.1182/blood-2019-123633>
13. Deng MC. The AlloMap™ genomic biomarker story: 10 years after. *Clin Transplant.* 2017;31(3).
14. He J, Abdel-Wahab O, Nahas MK, Wang K, Rampal RK, Intlekofer AM, et al. Integrated genomic DNA/RNA profiling of hematologic malignancies in the clinical setting. *Blood.* 2016;127(24):3004–14.
15. Ma XJ, Patel R, Wang X, Salunga R, Murage J, Desai R, et al. Molecular classification of human cancers using a 92-gene real-time quantitative polymerase chain reaction assay. *Arch Pathol Lab Med.* 2006;130(4):465–73.
16. McShane LM, Cavenagh MM, Lively TG, Eberhard DA, Bigbee WL, Williams PM, et al. Criteria for the use of omics-based predictors in clinical trials. *Nature.* 2013.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

17. Moons KGM, Altman DG, Reitsma JB, Ioannidis JPA, Macaskill P, Steyerberg EW, et al. Transparent reporting of a multivariable prediction model for individual prognosis or diagnosis (TRIPOD): Explanation and elaboration. *Ann Intern Med*. 2015;162(1):W1–73.
18. Banzi R, Gerardi C, Fratelli M, Garcia P, Torres T, Abad JMH, et al. Web-page for the Personalized Medicine Trials (PERMIT) project [Internet]. 2020 [cited 2021 Aug 2]. Available from: <https://permit-eu.org>
19. Banzi R, Gerardi C, Fratelli M, Garcia P, Torres T, Abad JMH, et al. Methodological approaches for personalised medicine: protocol for a series of scoping reviews [Internet]. 10.5281/zenodo.3770937. Available from: <https://zenodo.org/record/3770937>
20. Perlis RH. Translating biomarkers to clinical practice. Vol. 16, *Molecular Psychiatry*. 2011. p. 1076–87.
21. Graaf G, Postmus D, Westerink J, Buskens E. The early economic evaluation of novel biomarkers to accelerate their translation into clinical applications. *Cost Eff Resour Alloc*. 2018;16(1).
22. Drucker E, Krapfenbauer K. Pitfalls and limitations in translation from biomarker discovery to clinical utility in predictive and personalised medicine. Vol. 4, *EPMA Journal*. 2013.
23. Williams JK, Anderson CM. Omics research ethics considerations. *Nurs Outlook*. 2018;
24. Vähäkangas K. Research ethics in the post-genomic era. *Environmental and Molecular Mutagenesis*. 2013.
25. Olson RS, La Cava W, Orzechowski P, Urbanowicz RJ, Moore JH. PMLB: A large benchmark suite for machine learning evaluation and comparison. *BioData Min*. 2017;
26. Rappoport N, Shamir R. Multi-omic and multi-view clustering algorithms: Review and cancer benchmark. *Nucleic Acids Res* [Internet]. 2018;46(20):10546–62. Available from: <http://www.embase.com/search/results?subaction=viewrecord&from=export&id=L626271229>
27. Herrmann M, Probst P, Hornung R, Jurinovic V, Boulesteix AL. Large-scale benchmark study of survival prediction methods using multi-omics data. *Brief Bioinform*. 2021;
28. Peters MDJ, Godfrey CM, Khalil H, McInerney P, Parker D, Soares CB. Guidance for conducting systematic scoping reviews. *Int J Evid Based Healthc*. 2015;13(3):141–6.
29. Pham MT, Rajić A, Greig JD, Sargeant JM, Papadopoulos A, Mcewen SA. A scoping review of scoping reviews: Advancing the approach and enhancing the consistency. *Res Synth Methods*. 2014;5(4):371–85.
30. Munn Z, Peters MDJ, Stern C, Tufanaru C, McArthur A, Aromataris E. Systematic review or scoping review? Guidance for authors when choosing between a systematic or scoping review approach. *BMC Med Res Methodol*. 2018;18(1):143.
31. Tricco AC, Lillie E, Zarin W, O'Brien KK, Colquhoun H, Levac D, et al. PRISMA extension for scoping reviews (PRISMA-ScR): Checklist and explanation. *Ann Intern Med*. 2018 Oct;169(7):467–73.
32. Van't Veer LJ, Dai H, Van de Vijver MJ, He YD, Hart AAM, Mao M, et al. Gene expression profiling predicts clinical outcome of breast cancer. *Nature*. 2002;415(6871):530–6.
33. Ouzzani M, Hammady H, Fedorowicz Z, Elmagarmid A. Rayyan-a web and mobile app for systematic reviews. *Syst Rev*. 2016;5(1):210.
- [dataset] 34. Glaab E, Armin, Rita, Chiara, Paula, Jacques, et al. Data from: Selected articles from the scoping review of biomarker discovery studies for the EU project on "Personalised Medicine Trials" (PERMIT). Zenodo, November 4, 2021.

1
2
3 <https://doi.org/10.5281/zenodo.5646467>

- 4
5 35. Rahman M, Fukui T. Biomedical research productivity: Factors across the countries. *Int J Technol Assess Health Care*. 2003;19(1):249–52.
- 6
7 36. Hwang TJ, Kesselheim AS, Vokinger KN. Lifecycle Regulation of Artificial Intelligence- And Machine Learning-Based Software Devices in Medicine. Vol. 322, *JAMA - Journal of the American Medical Association*. 2019. p. 2285–6.
- 8
9
10
11 37. FDA Center for Devices and Radiological Health. Web-page on Nucleic Acid Based Tests by the Food and Drug Administration (FDA) [Internet]. 2021 [cited 2021 Aug 2]. Available from: <https://www.fda.gov/medical-devices/vitro-diagnostics/nucleic-acid-based-tests>
- 12
13
14
15 38. Winner BS, Sgroi DC, Ryan PD, Bruinsma TJ, Glas AM, Male A, et al. Analysis of the mamma print breast cancer assay in a predominantly postmenopausal cohort. *Clin Cancer Res*. 2008;14(10):2988–93.
- 16
17
18
19 39. Slodkowska EA, Ross JS. MammaPrint™ 70-gene signature: Another milestone in personalized medical care for breast cancer patients. *Expert Rev Mol Diagn*. 2009;9(5):417–22.
- 20
21
22
23 40. Sapino A, Roepman P, Linn SC, Snel MHJ, Delahaye LJMJ, Van Den Akker J, et al. MammaPrint molecular diagnostics on formalin-fixed, paraffin-embedded tissue. *J Mol Diagnostics*. 2014;16(2):190–7.
- 24
25
26
27 41. Mook S, Knauer M, Bueno-De-Mesquita JM, Retel VP, Wesseling J, Linn SC, et al. Metastatic potential of T1 breast cancer can be predicted by the 70-gene MammaPrint signature. *Ann Surg Oncol*. 2010;17(5):1406–13.
- 28
29
30
31 42. Maak M, Simon I, Nitsche U, Roepman P, Snel M, Glas AM, et al. Independent validation of a prognostic genomic signature (ColoPrint) for patients with stage II colon cancer. *Ann Surg*. 2013;257(6):1053–8.
- 32
33
34
35 43. Kopetz S, Tabernero J, Rosenberg R, Jiang Z, Moreno V, Bachleitner-Hofmann T, et al. Genomic Classifier ColoPrint Predicts Recurrence in Stage II Colorectal Cancer Patients More Accurately Than Clinical Factors. *Oncologist*. 2015;20(2):127–33.
- 36
37
38
39 44. Tan IB, Tan P. Genetics: An 18-gene signature (ColoPrint®) for colon cancer prognosis. *Nat Rev Clin Oncol*. 2011;8(3):131–3.
- 40
41
42
43 45. Rosenberg R, Maak M, Simon I, Nitsche U, Schuster T, Kuenzli B, et al. Independent validation of a prognostic genomic profile (ColoPrint) for stage II colon cancer (CC) patients. *J Clin Oncol*. 2011;29(4_suppl):358–358.
- 44
45
46
47 46. Salazar R, de Waard JW, Glimelius B, Marshall J, Klaase J, Van Der Hoeven J, et al. The PARSC trial, a prospective study for the assessment of recurrence risk in stage II colon cancer (CC) patients using ColoPrint. *J Clin Oncol*. 2012;30(4_suppl):678–678.
- 48
49
50
51 47. Tabernero J, Moreno V, Rosenberg R, Nitsche U, Bachleitner-Hofmann T, Lanza G, et al. Clinical and technical validation of a genomic classifier (ColoPrint) for predicting outcome of patients with stage II colon cancer. *J Clin Oncol*. 2012;30(4_suppl):384–384.
- 52
53
54
55 48. Bachleitner-Hofmann T, Simon I, Salazar R, Tabernero J, Rosenberg R, van der Akker J, et al. Development and Validation of a Robust Molecular Diagnostic Test (COLOPRINT) for Predicting Outcome in Stage II Colon Cancer Patients. *Ann Oncol*. 2012;23:ix179.
- 56
57
58
59 49. Nielsen T, Wallden B, Schaper C, Ferree S, Liu S, Gao D, et al. Analytical validation of the PAM50-based Prosigna Breast Cancer Prognostic Gene Signature Assay and nCounter Analysis System using formalin-fixed paraffin-embedded breast tumor specimens. *BMC Cancer*. 2014;14(1).
- 60 50. Wallden B, Storhoff J, Nielsen T, Dowidar N, Schaper C, Ferree S, et al. Development and

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

verification of the PAM50-based Prosigna breast cancer gene signature assay. *BMC Med Genomics*. 2015;8(1).

51. Alvarado MD, Prasad C, Rothney M, Cherbavaz DB, Sing AP, Baehner FL, et al. A Prospective Comparison of the 21-Gene Recurrence Score and the PAM50-Based Prosigna in Estrogen Receptor-Positive Early-Stage Breast Cancer. *Adv Ther*. 2015;32(12):1237–47.
52. Jensen MB, Lænkholm AV, Nielsen TO, Eriksen JO, Wehn P, Hood T, et al. The Prosigna gene expression assay and responsiveness to adjuvant cyclophosphamide-based chemotherapy in premenopausal high-risk patients with breast cancer. *Breast Cancer Res*. 2018;20(1).
53. Hequet D, Callens C, Gentien D, Albaud B, Mouret-Reynier MA, Dubot C, et al. Prospective, multicenter French study evaluating the clinical impact of the Breast Cancer Intrinsic Subtype-Prosigna® Test in the management of early-stage breast cancers. *PLoS One*. 2017;12(10).
54. Parker JS, Mullins M, Cheang MCU, Leung S, Voduc D, Vickery T, et al. Supervised risk predictor of breast cancer based on intrinsic subtypes. *J Clin Oncol*. 2009;27(8):1160–7.
55. Bartlett JMS, Bayani J, Marshall A, Dunn JA, Campbell A, Cunningham C, et al. Comparing Breast Cancer Multiparameter Tests in the OPTIMA Prelim Trial: No Test is More Equal Than the Others. *J Natl Cancer Inst*. 2016;108(9).
56. Kelly CM, Krishnamurthy S, Bianchini G, Litton JK, Gonzalez-Angulo AM, Hortobagyi GN, et al. Utility of oncotype DX risk estimates in clinically intermediate risk hormone receptor-positive, HER2-normal, grade II, lymph node-negative breast cancers. *Cancer*. 2010;116(22):5161–7.
57. Lo SS, Mumby PB, Norton J, Rychlik K, Smerage J, Kash J, et al. Prospective multicenter study of the impact of the 21-gene recurrence score assay on medical oncologist and patient adjuvant breast cancer treatment selection. *J Clin Oncol*. 2010;28(10):1671–6.
58. Carlson JJ, Roth JA. The impact of the Oncotype Dx breast cancer assay in clinical practice: A systematic review and meta-analysis. Vol. 141, *Breast Cancer Research and Treatment*. 2013. p. 13–22.
59. Thakur SS, Li H, Chan AMY, Tudor R, Bigras G, Morris D, et al. The use of automated Ki67 analysis to predict Oncotype DX risk-of-recurrence categories in early-stage breast cancer. *PLoS One*. 2018/01/06. 2018;13(1):e0188983.
60. Pease AM, Riba LA, Gruner RA, Tung NM, James TA. Oncotype DX® Recurrence Score as a Predictor of Response to Neoadjuvant Chemotherapy. *Ann Surg Oncol*. 2019;
61. Gianni L, Zambetti M, Clark K, Baker J, Cronin M, Wu J, et al. Gene expression profiles in paraffin-embedded core biopsy tissue predict response to chemotherapy in women with locally advanced breast cancer. *J Clin Oncol*. 2005;23(29):7265–77.
62. Paik S, Shak S, Tang G, Kim C, Baker J, Cronin M, et al. A multigene assay to predict recurrence of tamoxifen-treated, node-negative breast cancer. *N Engl J Med*. 2004;351(27):2817–26.
63. Dowsett M, Sestak I, Lopez-Knowles E, Sidhu K, Dunbier AK, Cowens JW, et al. Comparison of PAM50 risk of recurrence score with oncotype DX and IHC4 for predicting risk of distant recurrence after endocrine therapy. *J Clin Oncol*. 2013;31(22):2783–90.
64. Marrone M, Potosky AL, Penson D, Freedman AN. A 22 gene-expression assay, decipher® (GenomeDx biosciences) to predict five-year risk of metastatic prostate cancer in men treated with radical prostatectomy. *PLoS Curr*. 2015;7(EVIDENCEONGENOMICTESTS).
65. Nguyen PL, Haddad Z, Lam LLC, Ong K, Buerki C, Deheshi S, et al. Evaluation of the

- 1
2
3 Decipher prostate cancer classifier to predict metastasis and disease-specific mortality from
4 genomic analysis of diagnostic prostate needle biopsy specimens. *J Clin Oncol*.
5 2017;35(6_suppl):4–4.
6
- 7 66. Magi-Galluzzi C, Yousefi K, Haddad Z, Palmer-Aronsten B, Lam LLC, Buerki C, et al.
8 Validation of the Decipher prostate cancer classifier for predicting 10-year postoperative
9 metastasis from analysis of diagnostic needle biopsy specimens. *J Clin Oncol*.
10 2016;34(2_suppl):59–59.
11
- 12 67. Dalela D, Löppenber B, Sood A, Sammon J, Abdollah F. Contemporary Role of the
13 Decipher® Test in Prostate Cancer Management: Current Practice and Future Perspectives.
14 *Rev Urol*. 2016;18(1):1–9.
15
- 16 68. Klein EA, Haddad Z, Yousefi K, Lam LLC, Wang Q, Choeurng V, et al. Decipher Genomic
17 Classifier Measured on Prostate Biopsy Predicts Metastasis Risk. *Urology*. 2016;90:148–52.
18
- 19 69. Weiss LM, Chu P, Schroeder BE, Singh V, Zhang Y, Erlander MG, et al. Blinded comparator
20 study of immunohistochemical analysis versus a 92-gene cancer classifier in the diagnosis
21 of the primary site in metastatic tumors. *J Mol Diagnostics*. 2013;15(2):263–9.
22
- 23 70. Greco FA, Spigel DR, Yardley DA, Erlander MG, Ma X, Hainsworth JD. Molecular Profiling
24 in Unknown Primary Cancer: Accuracy of Tissue of Origin Prediction. *Oncologist*.
25 2010;15(5):500–6.
26
- 27 71. Hainsworth JD, Rubin MS, Spigel DR, Boccia R V., Raby S, Quinn R, et al. Molecular gene
28 expression profiling to predict the tissue of origin and direct site-specific therapy in patients
29 with carcinoma of unknown primary site: A prospective trial of the Sarah cannon research
30 institute. *J Clin Oncol*. 2013;31(2):217–23.
31
- 32 72. Harrison G, Sosa JA, Jiang X. Evaluation of the Afirma gene expression classifier in repeat
33 indeterminate thyroid nodules. *Arch Pathol Lab Med*. 2017;141(7):985–9.
34
- 35 73. Chudova D, Wilde JI, Wang ET, Wang H, Rabbee N, Egidio CM, et al. Molecular
36 classification of thyroid nodules using high-dimensionality genomic data. *J Clin Endocrinol
37 Metab*. 2010;95(12):5296–304.
38
- 39 74. Kim MI, Alexander EK. Diagnostic use of molecular markers in the evaluation of thyroid
40 nodules. Vol. 18, *Endocrine Practice*. 2012. p. 796–802.
41
- 42 75. Ali SZ, Fish SA, Lanman R, Randolph GW, Sosa JA. Use of the Afirma® gene expression
43 classifier for preoperative identification of benign thyroid nodules with indeterminate fine
44 needle aspiration cytology. *PLoS Currents*. 2013. p. 1–7.
45
- 46 76. McIver B, Castro MR, Morris JC, Bernet V, Smallridge R, Henry M, et al. An independent
47 study of a gene expression classifier (Afirma) in the evaluation of cytologically indeterminate
48 thyroid nodules. *J Clin Endocrinol Metab*. 2014;99(11):4069–77.
49
- 50 77. Lastra RR, Pramick MR, Crammer CJ, LiVolsi VA, Baloch ZW. Implications of a suspicious
51 Afirma test result in thyroid fine-needle aspiration cytology: An institutional experience.
52 *Cancer Cytopathol*. 2014;122(10):737–44.
53
- 54 78. Kim JY, Park SC, Lee JK, Choi SJ, Hahm KS, Park Y. Novel Antibacterial Activity of β 2-
55 Microglobulin in Human Amniotic Fluid. *PLoS One*. 2012;
56
- 57 79. Frampton GM, Fichtenholtz A, Otto GA, Wang K, Downing SR, He J, et al. Development
58 and validation of a clinical cancer genomic profiling test based on massively parallel DNA
59 sequencing. *Nat Biotechnol*. 2013;31(11):1023–31.
60
80. Wang K, Sanchez-Martin M, Wang X, Knapp KM, Koche R, Vu L, et al. Patient-derived
xenotransplants can recapitulate the genetic driver landscape of acute leukemias.
Leukemia. 2017;31(1):151–8.

- 1
2
3 81. Tarlock K, He J, Zhong S, Ries RE, Bailey M, Morley S, et al. Distinct age-associated
4 genomic profiles in acute myeloid leukemia (AML) using FoundationOne heme. *J Clin*
5 *Oncol*. 2016;34(15_suppl):7041–7041.
6
- 7 82. Lieber DS, Kennedy MR, Johnson DB, Greenbowe JR, Frampton GM, Schrock AB, et al.
8 Abstract B16: Validation and clinical feasibility of a Foundation Medicine assay to identify
9 immunotherapy response potential through tumor mutational burden (TMB). In 2017.
10
- 11 83. Lee Deak K, Jackson JB, Valkenburg KC, Keefer LA, Robinson Gerding KM, Angiuoli SV, et
12 al. Next-Generation Sequencing Concordance Analysis of Comprehensive Solid Tumor
13 Profiling between a Centralized Specialty Laboratory and the Decentralized PGDx Elio
14 Tissue Complete Kitted solution. *J Mol Diagnostics* [Internet]. 2021 Jul;in press. Available
15 from: <https://linkinghub.elsevier.com/retrieve/pii/S1525157821002105>
16
- 17 84. Labriola MK, Zhu J, Gupta R, McCall S, Jackson J, Kong EF, et al. Characterization of
18 tumor mutation burden, PD-L1 and DNA repair genes to assess relationship to immune
19 checkpoint inhibitors response in metastatic renal cell carcinoma. *J Immunother Cancer*
20 [Internet]. 2020 Mar;8(1):e000319. Available from:
21 <https://jitc.bmj.com/lookup/doi/10.1136/jitc-2019-000319>
22
- 23 85. Yamani MH, Taylor DO, Rodriguez ER, Cook DJ, Zhou L, Smedira N, et al. Transplant
24 Vasculopathy Is Associated With Increased AlloMap Gene Expression Score. *J Hear Lung*
25 *Transplant*. 2007;26(4):403–6.
26
- 27 86. Yamani MH, Taylor DO, Haire C, Smedira N, Starling RC. Post-transplant ischemic injury is
28 associated with up-regulated AlloMap gene expression. *Clin Transplant*. 2007;21(4):523–5.
29
- 30 87. Kobashigawa J, Patel J, Azarbal B, Kittleson M, Chang D, Czer L, et al. Randomized Pilot
31 Trial of Gene Expression Profiling Versus Heart Biopsy in the First Year after Heart
32 Transplant: Early Invasive Monitoring Attenuation Through Gene Expression Trial. *Circ Hear*
33 *Fail*. 2015;8(3):557–64.
34
- 35 88. Wingrove JA, Daniels SE, Sehnert AJ, Tingley W, Elashoff MR, Rosenberg S, et al.
36 Correlation of peripheral-blood gene expression with the extent of coronary artery stenosis.
37 *Circ Cardiovasc Genet*. 2008;1(1):31–8.
38
- 39 89. Rosenberg S, Dehais C, Ducray F, Alentron A, Tanguy M, De Reyneis A, et al. OS11.3
40 Machine learning for better prognostic stratification and driver genes identification in 1p/19q-
41 codeleted grade III gliomas. *Neuro Oncol* [Internet]. 2017;19(suppl_3):iii22–iii22. Available
42 from: <http://dx.doi.org/10.1093/neuonc/nox036>
43
- 44 90. Vargas J, Lima JAC, Kraus WE, Douglas PS, Rosenberg S. Use of the Corus® CAD Gene
45 Expression Test for Assessment of Obstructive Coronary Artery Disease Likelihood in
46 Symptomatic Non-Diabetic Patients. *PLoS Currents*. 2013.
47
- 48 91. Elashoff MR, Wingrove JA, Beineke P, Daniels SE, Tingley WG, Rosenberg S, et al.
49 Development of a blood-based gene expression algorithm for assessment of obstructive
50 coronary artery disease in non-diabetic patients. *BMC Med Genomics*. 2011;4.
51
- 52 92. Rosenberg S, Elashoff MR, Beineke P, Daniels SE, Wingrove JA, Tingley WG, et al.
53 Multicenter validation of the diagnostic accuracy of a blood-based gene expression test for
54 assessing obstructive coronary artery disease in nondiabetic patients. *Ann Intern Med*.
55 2010;153(7):425–34.
56
- 57 93. Brahe CH, Østergaard M, Johansen JS, Defranoux N, Wang X, Bolce R, et al. Predictive
58 value of a multi-biomarker disease activity score for clinical remission and radiographic
59 progression in patients with early rheumatoid arthritis: a post-hoc study of the OPERA trial.
60 *Scand J Rheumatol*. 2019;
94. Chernoff D, Scott Eastman P, Hwang CC, Flake DD, Wang X, Kivitz A, et al. Determination
of the minimally important difference (MID) in multi-biomarker disease activity (MBDA) test

- 1
2
3 scores: impact of diurnal and daily biomarker variation patterns on MBDA scores. *Clin*
4 *Rheumatol.* 2019;38(2):437–45.
5
6 95. Curtis JR, Weinblatt ME, Shadick NA, Brahe CH, Østergaard M, Hetland ML, et al.
7 Validation of the adjusted multi-biomarker disease activity score as a prognostic test for
8 radiographic progression in rheumatoid arthritis: a combined analysis of multiple studies.
9 *Arthritis Res Ther.* 2021;
10
11 96. Curtis JR, Xie F, Yang S, Danila MI, Owensby JK, Chen L. Uptake and Clinical Utility of
12 Multibiomarker Disease Activity Testing in the United States. *J Rheumatol* [Internet].
13 2019;46(3):237–44. Available from: <http://www.ncbi.nlm.nih.gov/pubmed/30442830>
14
15 97. Curtis JR, Van Der Helm-Van Mil AH, Knevel R, Huizinga TW, Haney DJ, Shen Y, et al.
16 Validation of a novel multibiomarker test to assess rheumatoid arthritis disease activity.
17 *Arthritis Care Res.* 2012;64(12):1794–803.
18
19 98. Food and Drug Administration. Helix Genetic Health Risk App For Late-Onset Alzheimer's
20 Disease - FDA Review Decision Summary. 2020; Available from:
21 <https://www.accessdata.fda.gov/scripts/cdrh/cfdocs/cfpmn/pmn.cfm?ID=K192073>
22
23 99. Cirulli ET, White S, Read RW, Elhanan G, Metcalf WJ, Tanudjaja F, et al. Genome-wide
24 rare variant analysis for thousands of phenotypes in over 70,000 exomes from two cohorts.
25 *Nat Commun.* 2020;
26
27 100. Lu JT, Ferber M, Hagenkord J, Levin E, South S, Kang HP, et al. Evaluation for Genetic
28 Disorders in the Absence of a Clinical Indication for Testing: Elective Genomic Testing.
29 *Journal of Molecular Diagnostics.* 2019.
30
31 101. Grzymalski JJ, Elhanan G, Morales Rosado JA, Smith E, Schlauch KA, Read R, et al.
32 Population genetic screening efficiently identifies carriers of autosomal dominant diseases.
33 *Nat Med.* 2020;
34
35 102. Huang S, Chaudhary K, Garmire LX. More is better: Recent progress in multi-omics data
36 integration methods. Vol. 8, *Frontiers in Genetics.* 2017.
37
38 103. Bersanelli M, Mosca E, Remondini D, Giampieri E, Sala C, Castellani G, et al. Methods for
39 the integration of multi-omics data: mathematical aspects. *BMC Bioinformatics* [Internet].
40 2016;17(S2):S15. Available from:
41 <http://bmcbioinformatics.biomedcentral.com/articles/10.1186/s12859-015-0857-9>
42
43 104. Picard M, Scott-Boyer MP, Bodein A, Périn O, Droit A. Integration strategies of multi-omics
44 data for machine learning analysis. Vol. 19, *Computational and Structural Biotechnology*
45 *Journal.* 2021. p. 3735–46.
46
47 105. Lemke C, Budka M, Gabrys B. Metalearning: a survey of trends and technologies. *Artif Intell*
48 *Rev.* 2015;
49
50 106. Sima C, Braga-Neto UM, Dougherty ER. High-dimensional bolstered error estimation.
51 *Bioinformatics.* 2011;
52
53 107. Ali M, Aittokallio T. Machine learning and feature selection for drug response prediction in
54 precision oncology applications. *Biophysical Reviews.* 2019.
55
56 108. Barla A, Jurman G, Riccadonna S, Merler S, Chierici M, Furlanello C. Machine learning
57 methods for predictive proteomics. *Br Bioinform.* 2008/03/04. 2008;9(2):119–28.
58
59 109. Grollemund V, Pradat PF, Querin G, Delbot F, Le Chat G, Pradat-Peyre JF, et al. Machine
60 learning in amyotrophic lateral sclerosis: Achievements, pitfalls, and future directions. *Front*
Neurosci [Internet]. 2019;13. Available from: <http://dx.doi.org/10.3389/fnins.2019.00135>
110. Ho DSW, Schierding W, Wake M, Saffery R, O'Sullivan J. Machine learning SNP based
prediction for precision medicine. *Front Genet* [Internet]. 2019;10(MAR). Available from:

- 1
2
3 <http://dx.doi.org/10.3389/fgene.2019.00267>
4
5 111. Krittanawong C, Bomback AS, Baber U, Bangalore S, Messerli FH, Wilson Tang WH. Future
6 Direction for Using Artificial Intelligence to Predict and Manage Hypertension. *Curr*
7 *Hypertens Rep*. 2018/07/08. 2018;20(9):75.
8
9 112. Long NP, Yoon SJ, Anh NH, Nghi TD, Lim DK, Hong YJ, et al. A systematic review on
10 metabolomics-based diagnostic biomarker discovery and validation in pancreatic cancer.
11 *Metabolomics*. 2019/03/05. 2018;14(8):109.
12
13 113. Martinez BI, Stabenfeldt SE. Current trends in biomarker discovery and analysis tools for
14 traumatic brain injury. *J Biol Eng [Internet]*. 2019;13(1). Available from:
15 <http://dx.doi.org/10.1186/s13036-019-0145-8>
16
17 114. Patil S, Awan KH, Arakeri G, Seneviratne CJ, Muddur N, Malik S, et al. Machine learning
18 and its potential applications to the genomic study of head and neck cancer-A systematic
19 review. *J Oral Pathol Med*. 2019;48(9):773–9.
20
21 115. Saini G, Mittal K, Rida P, Janssen EAM, Gogineni K, Aneja R. Panoptic view of prognostic
22 models for personalized breast cancer management. *Cancers (Basel) [Internet]*. 2019;11(9).
23 Available from: <http://dx.doi.org/10.3390/cancers11091325>
24
25 116. European Council. Council conclusions on personalised medicine for patients. *Off J Eur*
26 *Union [Internet]*. 2015;431(2):1–4. Available from: [https://eur-lex.europa.eu/legal-](https://eur-lex.europa.eu/legal-content/EN/TXT/PDF/?uri=CELEX:52015XG1217(01)&from=EN)
27 [content/EN/TXT/PDF/?uri=CELEX:52015XG1217\(01\)&from=EN](https://eur-lex.europa.eu/legal-content/EN/TXT/PDF/?uri=CELEX:52015XG1217(01)&from=EN)
28
29 117. Labriola MK, Zhu J, Gupta R, McCall S, Jackson J, Kong EF, et al. Characterization of
30 tumor mutation burden, PD-L1 and DNA repair genes to assess relationship to immune
31 checkpoint inhibitors response in metastatic renal cell carcinoma. *J Immunother Cancer*.
32 2020;8(1).
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Author Contributions

Study conception and design: EG, AR.

Methodology: CG, RB.

Data collection and analysis: EG, AR

Original draft preparation: EG

Review and editing: AR, EG, PG, CG, JDM, RB.

Project supervision: PG

Funding acquisition: JDM.

All authors have read and revised the manuscript and approved the final version.

The members of the PERMIT group were involved in the preparation or revision of the joint protocol of the four scoping reviews of the PERMIT series, attended the joint workshop (consultation exercise) and are co-authors of the other scoping reviews of the PERMIT series.

Collaborators

PERMIT group:

1. Antonio L. Andreu
2. Florence Bietrix,
3. Florie Brion Bouvier
4. Montserrat Carmona Rodriguez
5. Maria del Mar Polo-de Santos,
6. Maddalena Fratelli,
7. Rainer Girgenrath,
8. Alexander Grundmann,
9. Josep Maria Haro,
10. Frank Hulstaert,
11. Iñaki Imaz-Iglesia,
12. Setefilla Luengo Matos
13. Emmet McCormack,
14. Albert Sanchez Niubo,
15. Emanuela Oldoni,
16. Raphael Porcher,
17. Vibeke Fosse,
18. Luis M. Sánchez-Gómez,
19. Lorena San Miguel,
20. Cecilia Superchi,
21. Teresa Torres,
22. Anna Monistrol Mula

Funding statement

This project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No. 874825.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Competing interests statement

None declared

Ethics approval

This study was based entirely on a scoping review of relevant published literature and did not require an ethics approval.

Patient consent

This study did not require consent from patients, because it does not use individual data.

Permission to reproduce material from other sources

This study has cited all references which are published and publicly available.

Data sharing statement

The study protocol was published on the online platform Zenodo (19). Copies of searches and data extraction sheets will be made publicly available on Zenodo as part of the database collection for all scoping reviews conducted in the PERMIT project.

Figure legends

Fig. 1. Keyword based search strategy for the scoping review. Four categories of keywords were defined to retrieve relevant articles from the biomedical literature on machine learning analyses of omics data for personalised medicine, which include a validation study (highlighted by the coloured boxes in the centre). For each category relevant keywords were determined, including controlled vocabulary terms from the Medical Subject Headings (MeSH) thesaurus by the US National Library of Medicine (upper and lower boxes with frames coloured according to the corresponding category). As indicated by the keyword “AND” in the centre, a conjunctive search was conducted, i.e., every retrieved article had to contain at least one keyword from each category. This strategy was adapted for searching the other databases.

Fig. 2. Study selection flow diagram. Flow diagram of the procedure for the scoping review article identification, screening, eligibility assessment, and final inclusion, according to the PRISMA scheme (31). Reasons for excluding full-text were not mutually exclusive.

Fig. 3. Validation methods used in omics biomarker studies. Stacked bar chart of the number of articles retrieved in the scoping review for different categories of validation methods used in the underlying biomarker studies (covering time periods from 2000 to 2021). The majority of studies use only internal cohort validation approaches, such as cross-validation (CV), training/test set split validation, resampling/bootstrapping-based validation, out-of-bag validation (for tree-based classifiers), and combinations of CV and test set validation within the same cohort. Studies with an external validation on an independent patient cohort (with or without an additional internal cross-validation) are still underrepresented, even in more recent time periods. All filtered full-text articles derived from the scoping review except for review articles were included in the analysis.

Fig. 4. Map representation of country statistics for the selected articles. The number of articles originating from different countries among the studies selected in the full-text review are visualized on a world map representation using a colour gradient from blue (1 article) to red (98 articles = maximum contribution by a single country; using a logarithmic colour gradient scale to highlight differences over a broad value range).

Fig. 5. Representation of study types among the selected articles. The percentage of articles describing case-control studies, therapy/drug response studies, differential diagnosis studies, prognostic and survival prediction studies, as well as review studies and other study types is represented as a pie chart.

Fig. 6. Characteristics of successful omics-based studies. Six main categories of design and implementation aspects that characterize successful omics-based biomarker development studies were identified (starting from the centre left in the figure and proceeding clockwise): 1) Adequacy of the study design & sample size selection; 2) Rigor and robustness of the statistical evaluation; 3) Clarity of scope and goals; 4) Completeness and reproducibility of the study documentation; 5) Interpretability and biological plausibility of the created predictive models; 6) Integration of prior biological knowledge into the model building and validation procedures.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Fig. 7. Recommended generic workflow for biomarker development using machine learning analysis of omics data. The machine learning analysis of omics data for biomarker discovery and validation should ideally involve dedicated quality control and pre-processing analyses, a dimension reduction using unsupervised feature selection (e.g. a variance filtering) or data transformation approaches (e.g. using a Principal Components Analysis), a cross-validation on the discovery cohort, and an external validation on a distinct validation cohort.

For peer review only

Tables

Name	Test approval type	Purpose (Data type used for discovery)	References
MammaPrint®	FDA-cleared Assay	breast cancer risk-of-recurrence assessment (DNA microarray gene expression data)	(6,38–41)
ColoPrint®	LDT	colon cancer development of distant metastasis prediction (DNA microarray gene expression data)	(42–47)
Prosigna® Assay / PAM50	FDA-cleared Assay	breast cancer risk of distant recurrence prediction (DNA microarray gene expression data)	(49–53)
Oncotype DX®	LDT	breast cancer risk-of-recurrence assessment (DNA microarray gene expression data)	(8,56–59)
Decipher®	LDT	prostate cancer metastatic risk prediction (DNA microarray gene expression data)	(9,64–68)
Cancer Type ID®	LDT	predict tumour type for cancers of unknown / uncertain diagnosis (DNA microarray gene expression data)	(15,69–71)
Afirma™ Gene Expression Classifier	LDT	discriminate between benign and cancerous thyroid nodules (DNA microarray gene expression data)	(72–77)
Foundation One® Heme	LDT	test for hematologic malignancies, sarcomas, or solid tumours (RNA and DNA sequencing data)	(14,79–81)
PGDx Elio™ Tissue Complete	FDA-cleared Assay	test to assess somatic mutations and tumour mutation burden for solid tumours (DNA sequencing data)	(83,117)
AlloMap® Heart	FDA-cleared Assay	identifying heart transplant recipients with risk of cellular rejection (DNA microarray gene expression data)	(13,85–87)
Corus® CAD	LDT	identify obstructive coronary artery disease (DNA microarray gene expression data)	(11,88–91)
Vectra® DA	LDT	multi-biomarker blood test for rheumatoid arthritis (immunoassay + clinical data, 396 candidate biomarkers derived from integrative data analysis)	(93–96)
Helix® Laboratory Platform & Health Risk App for Late-onset Alzheimer's	FDA-cleared medical device	whole exome sequencing constituent device based for reporting and interpreting general genetic health risks (DNA sequencing data)	(99–101)

Tab. 1. Examples of clinically approved omics-derived diagnostic or prognostic tests designs applied

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

to personalised medicine (synonyms for the same test are separated by the “/”-symbol). FDA-approval status was checked on the web-site by the FDA (37) and reflects the status as of July 2021.

For peer review only

“stratified medicine” OR biomarker* OR “precision medicine”
OR “personalized medicine” OR “personalised medicine”
OR “individualized Medicine” OR “individualised Medicine”
OR “individualized therapy” OR “individualised therapy” OR
“patient stratification” OR pharmacogenetics OR “patient
specific modeling” OR “personalized clinical decision
making” OR “personalised clinical decision making” OR
“prediction of response” OR “prediction of responses” OR
“Biomarkers”[Mesh] OR “Precision Medicine”[Mesh]

Genomics”[Mesh]) OR “Metabolomics”[Mesh]) OR
“Epigenomics”[Mesh]) OR “Microarray Analysis”[Mesh]) OR
“Mass Spectrometry”[Mesh] OR Omic* OR “omic based”
OR “multi omic” OR “multi omics” OR genomic* OR
transcriptomic* OR proteomic* OR metabolomic OR
lipidomic* OR epigenomic* OR microarray OR “RNA seq”
OR “mass spectrometry”)

PERSONALISED
MEDICINE

OMICS

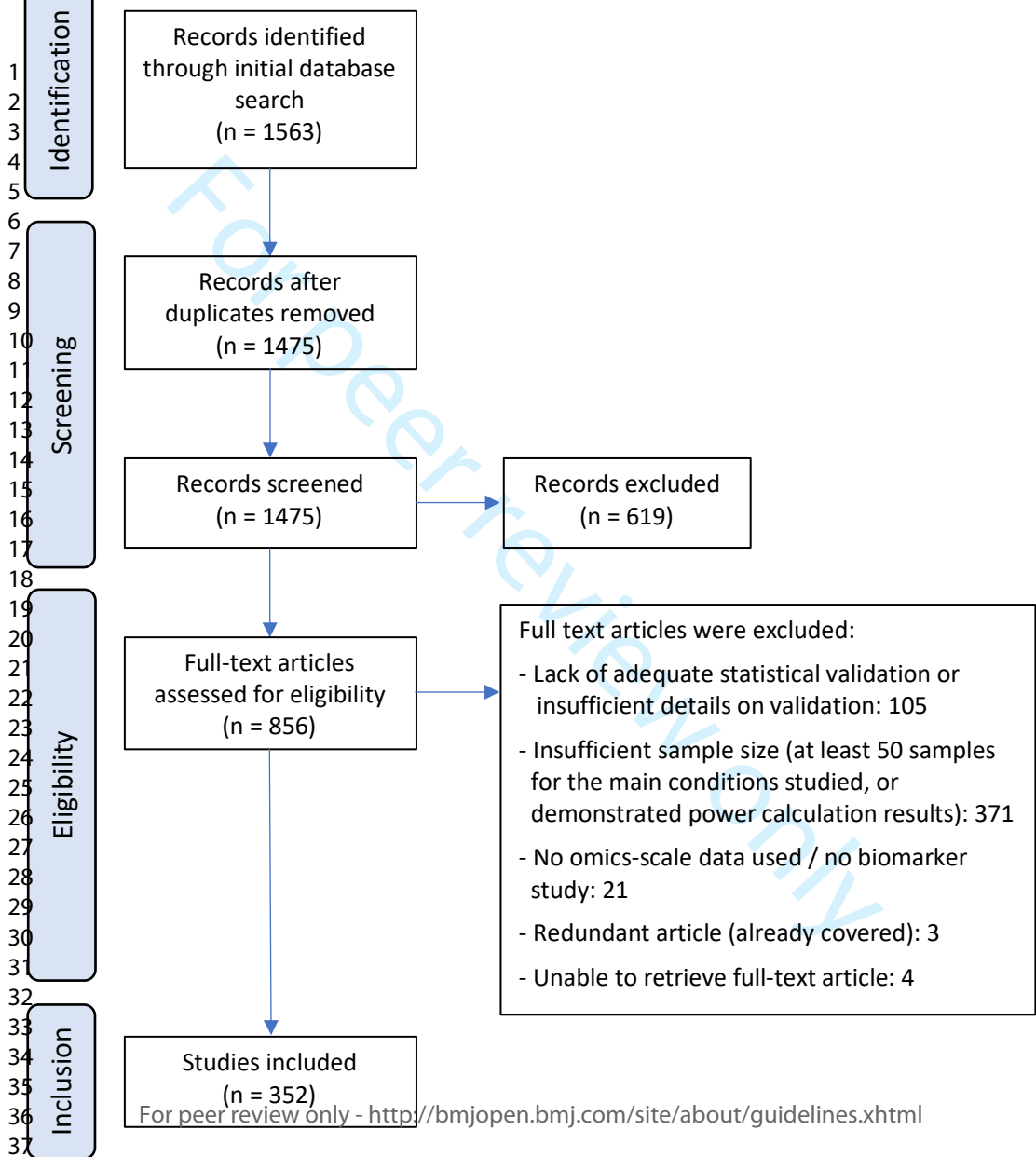
AND

MACHINE
LEARNING

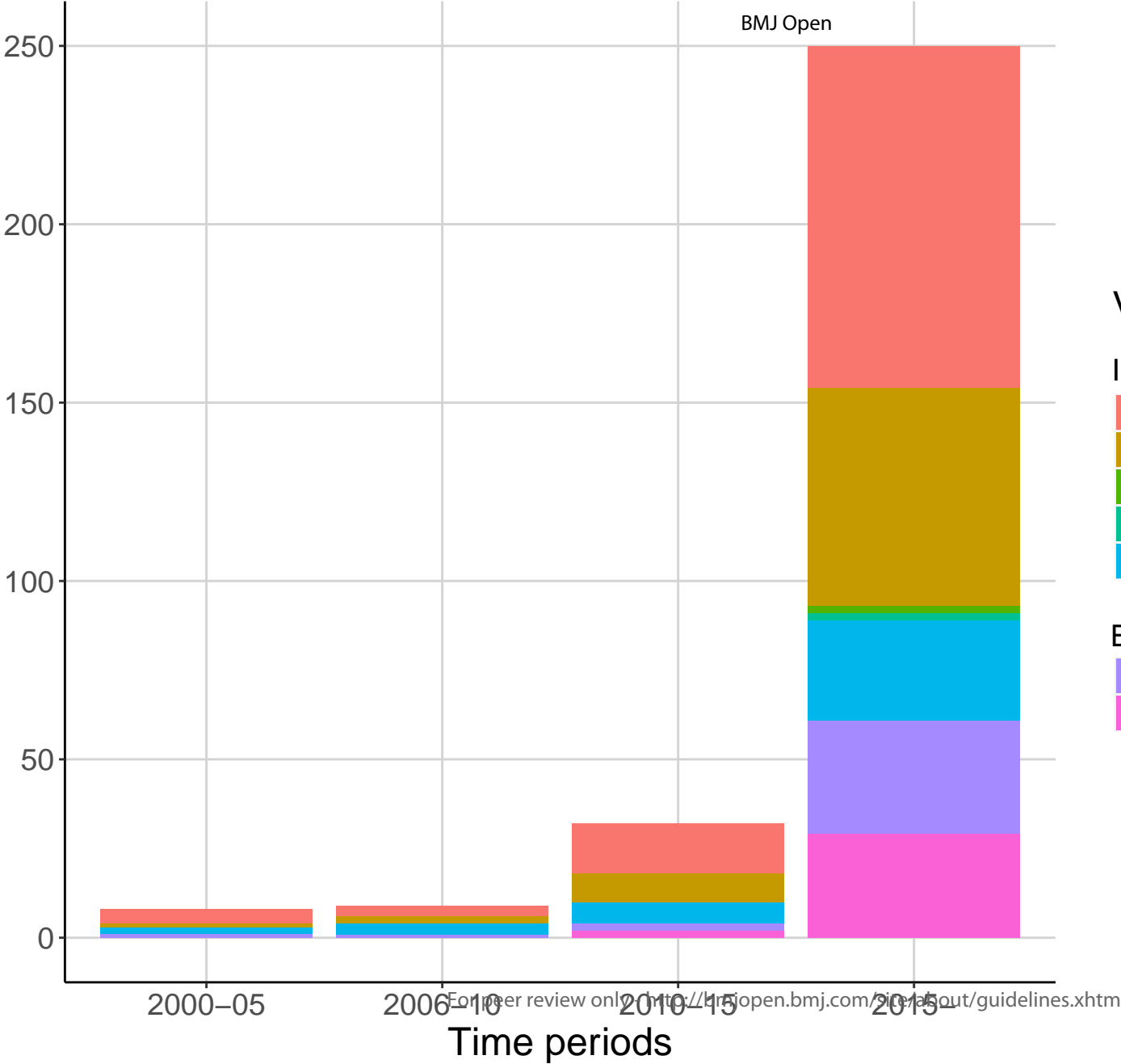
VALIDATION

Machine Learning”[Mesh] OR «Machine learning” OR
“statistical learning” OR “supervised learning” OR
“unsupervised learning”

Validation Studies as Topic”[Mesh]) OR “Validation Study”
[Publication Type] OR “Sensitivity and Specificity”[Mesh])
OR “Benchmarking”[Mesh]) OR validation OR validity OR
validated OR “cross validation” “cross validated” OR
“clinical utility*” OR accuracy OR robustness OR reliability*
OR sensitivity OR specificity OR benchmark* OR bias OR
“cross study” OR “cross studies”)



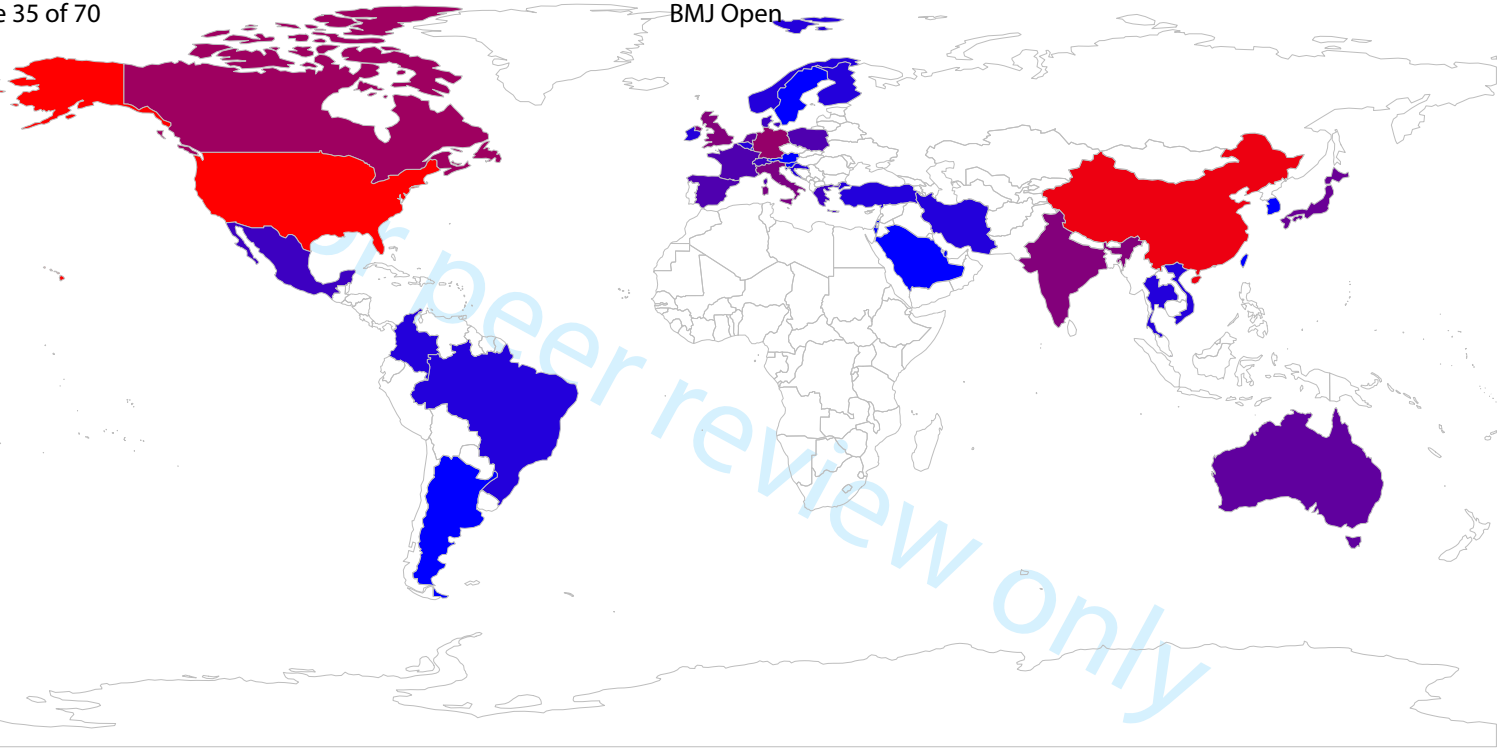
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38



- ### Validation methods
- Internal validation:
- Cross-validation (CV)
 - Training/test set validation
 - Resampling
 - Out-of-bag internal validation
 - CV internal cohort validation
- External validation:
- CV external cohort validation
 - External cohort validation

on 6 December 2021. Downloaded from <http://bmjopen.bmj.com/> on April 12, 2024 by guest. Protected by copyright.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26

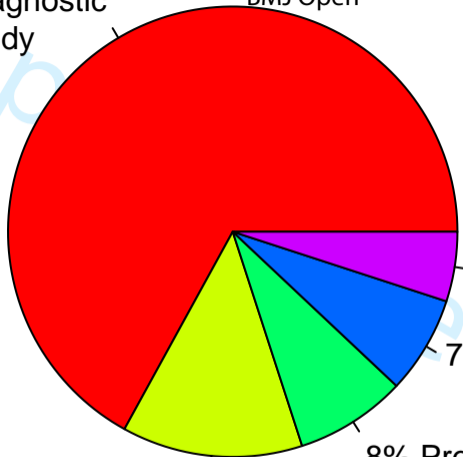


Downloaded from <http://bmjopen.bmj.com/> on April 10, 2024 by guest

67% Diagnostic study

BMJ Open

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16



5% Others

7% Therapy/drug response study

8% Prognostic/survival study

13% Review study

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29

For peer review only

Characteristics of successful omics-derived biomarker studies

Statistical evaluation

- cross-validation & external testing
- adequacy of performance & robustness metrics
- multiple testing correction

Clarity of scope & goals

- Inclusion/exclusion criteria
- primary/secondary outcomes

Study design & sample size

- statistical power
- balanced study groups
- batch effect avoidance/correction
- matching/adjustment for confounders/biases

Study documentation

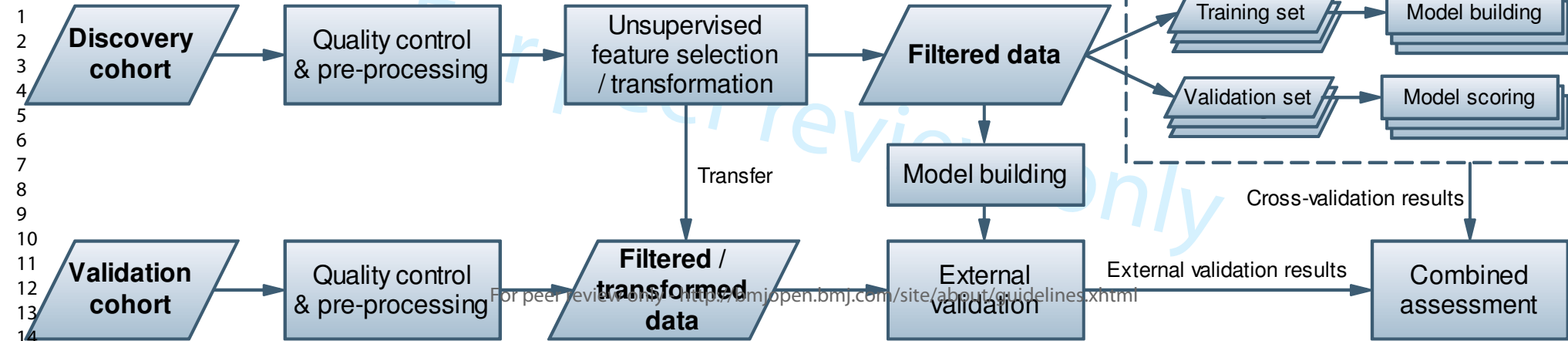
- instruments, settings & parameters
- reproducible methods description
- data provenance

Integration of prior knowledge

- molecular pathways & networks
- clinical data & real-world data
- multi-omics data

Model interpretability

- explainable predictions & human-interpretable models
- biological plausibility & mechanistic understanding



Online Supplementary file 1 – PRISMA-ScR Checklist

Preferred Reporting Items for Systematic reviews and Meta-Analyses extension for Scoping Reviews (PRISMA-ScR) Checklist (17).

SECTION	ITEM	PRISMA-ScR CHECKLIST ITEM	REPORTED ON PAGE #
TITLE			
Title	1	Identify the report as a scoping review.	1
ABSTRACT			
Structured summary	2	Provide a structured summary that includes (as applicable): background, objectives, eligibility criteria, sources of evidence, charting methods, results, and conclusions that relate to the review questions and objectives.	2
INTRODUCTION			
Rationale	3	Describe the rationale for the review in the context of what is already known. Explain why the review questions/objectives lend themselves to a scoping review approach.	4
Objectives	4	Provide an explicit statement of the questions and objectives being addressed with reference to their key elements (e.g., population or participants, concepts, and context) or other relevant key elements used to conceptualize the review questions and/or objectives.	4
METHODS			
Protocol and registration	5	Indicate whether a review protocol exists; state if and where it can be accessed (e.g., a Web address); and if available, provide registration information, including the registration number.	5
Eligibility criteria	6	Specify characteristics of the sources of evidence used as eligibility criteria (e.g., years considered, language, and publication status), and provide a rationale.	5-6
Information sources*	7	Describe all information sources in the search (e.g., databases with dates of coverage and contact with authors to identify additional sources), as well as the date the most recent search was executed.	5
Search	8	Present the full electronic search strategy for at least 1 database, including any limits used, such that it could be repeated.	5 (Online Suppl. File 2)
Selection of sources of evidence†	9	State the process for selecting sources of evidence (i.e., screening and eligibility) included in the scoping review.	6

SECTION	ITEM	PRISMA-ScR CHECKLIST ITEM	REPORTED ON PAGE #
Data charting process‡	10	Describe the methods of charting data from the included sources of evidence (e.g., calibrated forms or forms that have been tested by the team before their use, and whether data charting was done independently or in duplicate) and any processes for obtaining and confirming data from investigators.	6-7
Data items	11	List and define all variables for which data were sought and any assumptions and simplifications made.	6 (Online Suppl. File 3)
Critical appraisal of individual sources of evidence§	12	If done, provide a rationale for conducting a critical appraisal of included sources of evidence; describe the methods used and how this information was used in any data synthesis (if appropriate).	
Synthesis of results	13	Describe the methods of handling and summarizing the data that were charted.	6-7
RESULTS			
Selection of sources of evidence	14	Give numbers of sources of evidence screened, assessed for eligibility, and included in the review, with reasons for exclusions at each stage, ideally using a flow diagram.	7 (Fig. 2)
Characteristics of sources of evidence	15	For each source of evidence, present characteristics for which data were charted and provide the citations.	7 (Table 1)
Critical appraisal within sources of evidence	16	If done, present data on critical appraisal of included sources of evidence (see item 12).	
Results of individual sources of evidence	17	For each included source of evidence, present the relevant data that were charted that relate to the review questions and objectives.	7-13
Synthesis of results	18	Summarize and/or present the charting results as they relate to the review questions and objectives.	7-13
DISCUSSION			
Summary of evidence	19	Summarize the main results (including an overview of concepts, themes, and types of evidence available), link to the review questions and objectives, and consider the relevance to key groups.	13
Limitations	20	Discuss the limitations of the scoping review process.	14
Conclusions	21	Provide a general interpretation of the results with respect to the review questions and	14-16

SECTION	ITEM	PRISMA-ScR CHECKLIST ITEM	REPORTED ON PAGE #
		objectives, as well as potential implications and/or next steps.	
FUNDING			
Funding	22	Describe sources of funding for the included sources of evidence, as well as sources of funding for the scoping review. Describe the role of the funders of the scoping review.	25

For peer review only

Online Supplementary file 2 – Search strategy

Keyword searches conducted in the databases *PubMed*, *EMBASE* and *Web of Science* as part of the scoping review.

1) PubMed Query

Search: (((("Machine learning" OR "statistical learning" OR "supervised learning" OR "unsupervised learning") OR ("Machine Learning"[Mesh])) AND (("Biomarkers"[Mesh]) OR "Precision Medicine"[Mesh])) AND (((Omic* OR "omic based" OR "multi omic" OR "multi omics" OR genomic* OR transcriptomic* OR proteomic* OR metabolomic* OR lipidomic* OR epigenomic* OR microarray OR "RNA seq" OR "mass spectrometry")) OR ("Genomics"[Mesh] OR "Metabolomics"[Mesh] OR "Epigenomics"[Mesh] OR "Microarray Analysis"[Mesh] OR "Mass Spectrometry"[Mesh]))) AND ((validation OR validity OR validated OR "cross validation" "cross validated" OR "clinical utility*" OR accuracy OR robustness OR reliability* OR sensitivity OR specificity OR benchmark* OR bias OR "cross study" OR "cross studies")) AND ("2000/01/01"[Date - Entry]: "2021/07/20"[Date - Entry]) Filters: English, French, Italian, Spanish

2) Embase Query

#25: #24 AND [embase]/lim NOT [medline]/lim
 #24: #23 AND [2000-2021]/py
 #23: #20 AND #21 AND ([english]/lim OR [french]/lim OR [italian]/lim OR [spanish]/lim)
 #22: #20 AND #21
 #21: omic*:ti,ab OR 'machine learning':ti,ab OR 'personalized medicine':ti,ab OR 'personalised medicine':ti,ab
 #20: #4 AND #10 AND #16 AND #19
 #19: #17 OR #18
 #18: validation:ti,ab OR validity:ti,ab OR validated:ti,ab OR 'cross validation':ti,ab OR 'cross validated':ti,ab OR test*:ti,ab OR 'clinical utility*':ti,ab OR accuracy:ti,ab OR robustness:ti,ab OR reliability*:ti,ab OR sensitivity:ti,ab OR specificity:ti,ab OR benchmark*:ti,ab OR bias:ti,ab OR 'cross study:ti,ab' OR 'cross studies':ti,ab
 #17: 'validation study'/exp OR 'reliability'/exp OR 'sensitivity and specificity'/exp OR 'benchmarking'/exp
 #16: #14 OR #15
 #15: omic*:ti,ab OR 'omic based':ti,ab OR 'multi omic*':ti,ab OR genomic*:ti,ab OR transcriptomic*:ti,ab OR proteomic*:ti,ab OR metabolomic*:ti,ab OR lipidomic*:ti,ab OR epigenomic*:ti,ab OR microarray:ti,ab OR 'rna seq':ti,ab OR 'mass spectrometr*':ti,ab
 #14: #11 OR #12 OR #13
 #13: 'mass spectrometry'/exp
 #12: 'microarray analysis'/exp
 #11: 'omics'/exp OR 'genomics'/exp OR 'epigenetics'/exp
 #10: #5 OR #6 OR #7 OR #8 OR #9
 #9: 'individualized medicine':ti,ab OR 'individualised medicine':ti,ab OR 'individualized therapy':ti,ab OR 'individualised therapy':ti,ab
 #8: 'personalised medicine':ti,ab
 #7: 'personalized medicine':ti,ab
 #6: 'stratified medicine':ti,ab OR cluster*:ti,ab OR 'sub group*':ti,ab OR subgroup*:ti,ab OR biomarker*:ti,ab OR diagnos*:ti,ab OR prognos*:ti,ab OR 'precision medicine':ti,ab
 #5: 'biological marker'/exp OR 'personalized medicine'/exp
 #4: #1 OR #2 OR #3
 #3: 'machine learning'/exp
 #2: 'statistical learning'/exp

1
2
3 #1: 'machine learning':ti,ab OR 'statistical learning':ti,ab OR 'supervised learning':ti,ab OR
4 'unsupervised learning':ti,ab
5
6

7 **3) Web of Science Query**

8 (((((#1) AND #1) AND #2) AND #3) AND #4) AND #5

9 5: (ALL=(((validation OR validity OR validated OR "cross validation" "cross validated" OR
10 "clinical utility*" OR accuracy OR robustness OR reliability* OR sensitivity OR specificity OR
11 benchmark*

12 OR bias OR "cross study" OR "cross studies")))) AND ALL=(((omic* OR "machine learning"
13 OR

14 "personalized medicine" OR "personalised Medicine"))))

15 4: ALL=(((validation OR validity OR validated OR "cross validation" "cross validated" OR
16 "clinical utility*" OR accuracy OR robustness OR reliability* OR sensitivity OR specificity OR
17 benchmark*

18 OR bias OR "cross study" OR "cross studies"))))

19 3: ALL=(TOPIC: ((Omic* OR "omic based" OR "multi omic*" OR genomic* OR transcriptomic*
20 OR proteomic* OR metabolomic* OR lipidomic* OR epigenomic* OR microarray OR "RNA
21 seq"

22 OR "mass spectrometr*"))))

23 2: (ALL=(TOPIC: (("Machine learning" OR "statistical learning" OR "supervised learning" OR
24 "unsupervised learning")))) AND ALL=(TOPIC: (("stratified medicine" OR cluster* OR "sub

25 group*" OR Subgroup* OR biomarker* OR diagnos* OR prognos* OR "precision medicine"
26 OR

27 "personalized medicine"OR "personalised medicine" OR "individualized Medicine" OR
28 "individualised Medicine" OR "individualized therapy" OR "individualised therapy"))))

29 1: ALL=(TOPIC: (("Machine learning" OR "statistical learning" OR "supervised learning" OR
30 "unsupervised learning"))))
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Online Supplementary file 3 – Data extraction form

Data items extracted from each processed article during the full-text scoping review, and associated qualifications for each item.

Item	Qualifications
Authors	
Title	
Journal	
Volume	
Issue	(if applicable)
Pages	(if applicable)
Year	
Location	
URL / DOI	
Type of publication	<ul style="list-style-type: none"> • Research article • Meeting abstract • Review
Study population and sample size	(if applicable)
Methodology / Study Design	<ul style="list-style-type: none"> • Case-control study • Cases only stratification study <p>(+ further qualification, e.g. treatment response prediction, tumor subtype categorization, recurrence/relapse prediction, survival prediction, tissue-of-origin prediction)</p>
Outcome assessment	<ul style="list-style-type: none"> • Performance measures (e.g. accuracy, sensitivity, specificity, Kohen's Kappa, F-score, AUC) • Validation scheme (cross-validation approach, external validation approach, single cohort or multiple cohorts)
Generic machine learning category	<ul style="list-style-type: none"> • Supervised learning • Unsupervised learning • Other / mixed approaches
Name of specific machine learning approach	(if applicable)

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Main results / key findings that relate to the research question	short description
---	-------------------

For peer review only

1	Brown, K K and Choi, Y and Cully, T V and Flaherty, R and Grosberg, S and Imtiaz, U and Lynch, D A and Myers, J L and Steeg, P S and Martinez, F J and Parkin, D G and Wald, A J and Hwang, J and Barth, M M and Rago, G and Kennedy, G C	American Journal of Respiratory and Critical Care Medicine	195	2017	Prospective validation of a genomic classifier for usual interstitial pneumonia in transbronchial biopsies	meeting abstract	http://www.ajrccp.com/cgi/doi/10.1164/rccp.2017.195.11.1611	354 TB8 samples	Case-control study	AUC (training / test set split)	training + test set	
2												
3												
4												
5	25. Cai, Q and Alvarez, J A and Kang, J and Yu, T	J Proteome	16	3	1261-1269	2017	United States	Network Marker Selection for Untargeted LC-MS Metabolomics Data	Res	Cases only (BMI analysis)	AUC (5-fold CV)	cross-validation
6												
7												
8	26. Cai, Z and Xu, D and Zhang, Q and Zhang, J and Ngai, S M and Sha, J	Mol Biolyist	11	3	791-800	2015	China	Classification of lung cancer using ensemble-based feature selection and machine learning methods	Mol Biolyist	Case-control study	accuracy, precision, recall, F-score (LDOCV)	cross-validation
9												
10	Canouan, R and Varma, S and Simpson, B and Kim, M and An, Y and Saldana, S and Rivers, C and Moscato, P and Griswold, M and Sonntag, D and Waihele, J and Klavik, K and Jonsson, P V and Eriksson, G and Asplund, T and Luener, L J and Gudnason, V and Ungewill, C and Thambisetty, M	Alzheimers Dement	12	7	815-822	2016	Australia	Blood metabolite markers of preclinical Alzheimer's disease in two longitudinally followed cohorts of older individuals	Alzheimers Dement	Case-control study	AUC, sensitivity, specificity (5-fold CV)	cross-validation
11												
12												
13	28. Chai, Boonchoe, A and Samarasirige, S and Kulasi, D	Current Biomformats	5	2	118-133	2010		Machine Learning for Childhood Acute Lymphoblastic Leukemia Gene Expression Data Analysis: A Review	Current Biomformats	review (not applicable)		
14												
15	29. Chang, Y and Park, H and Yang, H J and Lee, S and Lee, K Y and Kim, T S and Jung, J and Shin, J M	Sci Rep	8	1	8857-8857	2018	Australia	Cancer Drug Response Profile scan (CDScan): A Deep Learning Model That Predicts Drug Effectiveness From Cancer Genomic Signature	Sci Rep	Cases only (drug response study)	Required, AUC (training/test split)	training + test set
16												
17												
18	30. Cho, S M and Connolly, J and Ng, Y L and Ganesan, I and Bennett, L	Pediatric Nephrology	31	10	1746-1746	2016		Can urinary proteomes be used as non-invasive markers for renal involvement in childhood fibrotic urinary tract infection (UTI)?	Pediatric Nephrology	Case-control study	sensitivity, PPV (10-fold CV)	cross-validation
19												
20	31. Chaudhary, A and Porion, D B and Li, L and Garmize, L X	Clin Cancer Res	24	6	1248-1259	2018	United States	Deep Learning-Based Multi-Omics Integration Robustly Predicts Survival in Liver Cancer	Clin Cancer Res	Cases only (survival prediction)		external cohort validation
21												
22	32. Choung, R S and Khaledi Rostamkobei, S and Ju, J M and Marietta, E V and Van Dyke, C T and Rajasekaran, J J and Jayaraman, V and Wang, T and Bei, K and Rajasekaran, K E and Krishna, K and Krishnamurthy, H K and Murray, J A	Gastroenterology	156	3	582-591.e1	2019	United States	Synthetic Neopeptides of the Transglutaminase-Deamidated Gliadin Complex as Biomarkers for Diagnosing and Monitoring Celiac Disease	Gastroenterology	Case-control study	AUC, accuracy, sensitivity, specificity ("We validated our findings in 82 patients with newly diagnosed celiac and 117 controls.")	external cohort validation
23												
24												
25												
26	33. Chang, W Y and Correa, E and Yoshimura, K and Chang, M C and Dennison, A and Takeda, S and Chang, Y T	American Journal of Translational Research	12	1	171-179	2020	Japan	Using probe electrospray ionization mass spectrometry and machine learning for detecting pancreatic cancer with high performance	American Journal of Translational Research	Case-control study	accuracy, sensitivity, specificity (1000 independent repetitions of a bootstrap cross-validation process)	cross-validation
27												
28												
29	34. Clark, O and Safikhani, Z and Srinivas, P and Kabe, Kian, B	Irish Journal of Medical Science	187		5348-5348	2018	Canada	Gene isoforms as expression-based biomarkers predictive of drug response in vitro	Irish Journal of Medical Science	Cases only (drug response prediction in vitro)	AUC, accuracy (validation in independent breast cancer data and different pharmacological assay)	external cohort validation
30												
31	35. Croner, J J and Kao, A and Benz, R and Blume, J E and Dillon, R and Wilcox, B and Kairs, S N	Clinical Chemistry	63		522-523	2017	Denmark	A new blood test for colorectal cancer in high-risk subjects	Clinical Chemistry	Case-control study	sensitivity, specificity, PPV, NPV (10-fold CV + training / test set split)	cross-validation + test set
32												
33												
34												
35	36. Cruz, J A and Wishart, D S	Cancer Informatics	2	2	59-77	2006	Greece	Applications of machine learning in cancer prediction and prognosis	Cancer Informatics	Case-control study	review (not applicable)	
36												
37												
38												
39												
40												
41	37. Cugliari, G and Benevenuto, S and Garretts, S and Sacardote, C and Panico, S and Krogh, V and Turmino, R and Virei, P and Fariello, P and Martello, G	Journal of Computational Intelligence in Informatics and Cybernetics	19	42	39-42	2019	USA	Improving the prediction of cancer risk with machine learning and DNA methylation data	Journal of Computational Intelligence in Informatics and Cybernetics	Case-control study	AUC, sensitivity, specificity (nested cross-validation)	cross-validation
42												
43												
44												
45												
46												

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46

38	Das, D and Ito, A and Kodawaki, T and Tsuda, K	An interpretable machine learning model for diagnosis of Alzheimer's disease	PeerJ	2019	3	2019	Japan	https://doi.org/10.7717/peerj.6641	article	97 AD subjects + 54 controls	Case-control study	AUC, accuracy, sensitivity, specificity (cross-validation + test set)	cross-validation + test set
39	de Maturana, E and Alonso, L and Alarcón, P and Martín-Antonián, I A and Pineda, S and Pirono, L and Calle, M I and Malats, N	Challenges in the integration of omics and non-omics data	Genes	10	3	2019	Spain	https://doi.org/10.3390/genes10030249	article	review (not applicable)	Review		
40	de Ronde, J J and Bonder, M J and Lips, E H and Rodenhuis, S and Wessels, L F	Breast cancer subtype specific classifiers of response to neoadjuvant chemotherapy do not outperform classifiers trained on all subtypes	PLoS One	9	2	2014		https://doi.org/10.1371/journal.pone.0085511	article	374 samples were analyzed	Cases only (treatment response prediction)	AUC (nested cross-validation)	cross-validation
41	Di Camillo, B and Sanavia, I and Martini, M and Jurman, G and Sambro, F and Baria, A and Squilarlo, M I and Furlanello, C and Toffolo, G and Cobelli, C	Effect of size and heterogeneity of samples on biomarker discovery: synthetic and real data assessment	PLoS One	7	3	2012	Italy	https://doi.org/10.1371/journal.pone.0032000	article	3 different datasets (more than 50 samples per group in total)	Case-control study	AUC, sensitivity, specificity (cross validation + Monte Carlo bootstrap resampling)	cross-validation
42	Diaz-Cano, S and Sutherland, R and Moorhead, J and Blanes, A and Dobson, R	Growth pattern analysis in low grade clear cell renal cell carcinoma: Prognostic value and biological significance	Laboratory Investigation	96		2016		https://doi.org/10.1038/liv.2016.08	meeting abstract	low FG (1-2, 174 cases) vs. high FG (3-4, 139 cases) grade	Subtype comparison	AUC (50-fold cross-validation)	cross-validation
43	Diggins, J and Kim, S Y and Hu, Z Z and Parkar, D and Wong, M and Reynolds, J and Tom, E and Pagan, M and Moore, R and Rossi, J and Livolsi, V A and Lamm, R B and Kloos, R T and Walsh, P S and Kennedy, G C	MACHINE LEARNING FROM CONCEPT TO CLINIC: RELIABLE DETECTION OF BRAF V600E DNA MUTATIONS IN THYROID NODULES USING HIGH-DIMENSIONAL RNA EXPRESSION DATA	Pacific Symposium on Biocomputing	2015			USA	https://doi.org/10.1146/annurev-bioinfor-050715-090077	article	training (n=181) and independent test (n=535)	Case-control study	AUC (10-fold CV + external test set)	cross-validation + test set
44	Ding, M Q and Chen, L and Cooper, G F and Young, J D and Liu, X	Precision Oncology beyond Targeted Therapy: Combining Omics Data with Machine Learning Matches the Specificity of Cancer Cells to Effective Therapeutics	Mol Cancer Res	16	2	2018	United States	https://doi.org/10.1007/s12013-018-02178-7	article	transcriptomics data from 277 cell lines was used	Cases only (cancer cell line drug response prediction)	accuracy, sensitivity, specificity (25-fold cross-validation)	cross-validation
45	Djebbari, A and Labbe, A	Refining gene signatures: a Bayesian approach	BMC Bioinformatics	10		2009	Canada	https://doi.org/10.1186/1471-2108-10-241	article	the approach was applied to multiple cancer microarray datasets with > 50 samples per group in total	Case-control study	AUC, sensitivity, specificity (10-fold CV + external test set)	cross-validation + test set
46	Dougherty, E R and Hua, J and Bittner, M L	Validation of computational methods in genomics	Current Genomics	8	1	2007	USA	https://doi.org/10.5555/0566	article	review (not applicable)	review		
47	Drouin, A and Giguère, S and Déry, M and Marchand, M and Yers, M and Lo, V G and Bourquart, A M and Laviolette, F and Corbeil, J	Predictive computational phenotyping and biomarker discovery using reference-free genome comparisons	BMC Genomics	17	1	2016	Canada	https://doi.org/10.1186/s12864-016-2886-6	article	17 datasets in which the number of examples ranged from 111 to 556	Antibiotic resistance prediction	error rate (5-fold CV, test evaluation)	cross-validation + test set
48	Drozdzal, I and Kidd, M and Modin, I	Graph-theoretic definition of neuroendocrine disease a tumor specific mathematical toolbox for assessing neoplastic behaviour	Neuroscience	103		2016		https://doi.org/10.1016/j.neuroscience.2016.04.048	meeting abstract	130 blood samples (NENs: n = 63)	Case-control study	AUC, sensitivity, specificity, PPV, NPV (The model was validated in two independent sets [Set 1 (n = 115, NENs: n = 72) Set 2 (n = 120, NENs: n = 58)])	training + test set
49	Ebeabali, E and Lee, F and Schendel, E and Houque, A and Kathirason, N and Pathare, T and Syed, N and Al-Ali, R	Large-scale machine learning based on functional networks for biomedical big data with high performance computing platforms	Journal of Computational Science	11		2015		https://doi.org/10.1016/j.jocs.2015.06.008	article	Stage1 (population of 127 exacerbation cases and 290 non-exacerbation controls) and Stage2 (population of 50 exacerbation cases and 114 non-exacerbation controls)	Case-control study	AUC, sensitivity, specificity (training/test set split)	training + test set
50	Fan, X and Wan, X B and Huang, Y and Cai, H M and Fu, X H and Yang, Z L and Chen, D K and Song, S X and Wu, P H and Liu, Q and Wang, L and Wang, J P	Epithelial-mesenchymal transition biomarkers and support vector machine guided model in preoperatively predicted rectal lymph node metastasis for rectal cancer	Br J Cancer	106	11	2012	China	https://doi.org/10.1038/bjc.2012.84	article	193 RC patients	Cases only (predicting lymph node metastasis)	accuracy, sensitivity, specificity (training/test set split)	training + test set
51	Fang, Y and Xu, P and Yang, J and Qin, Y	A quantile regression forest based method for drug response and assess prediction reliability	PLoS One	13	10	2018	China	https://doi.org/10.1371/journal.pone.0200110	article	data from 947 cell lines (CCLE dataset)	Cases only (drug response prediction in vitro)	Pearson correlation of observed and predicted drug response (out-of-bag prediction)	outofbag
52	Farmaki, D and Koek, T and Muller, W and Parisis, J and Gogas, B and Nikolaou, M and Lekakia, I and Mischak, H and Pippatos, G	Urine proteome analysis in heart failure with reduced ejection fraction complicated by chronic kidney disease: feasibility, and clinical and pathogenic correlates	Eur J Heart Fail	18	7	2016	Germany	https://doi.org/10.1002/ejhf.544	article	126 individuals, 59 HFpEF patients and 67 controls	Case-control study	AUC, accuracy, sensitivity, specificity (cross-validation + test set)	cross-validation + test set

"We present an interpretable machine learning model for medical diagnosis called sparse high-order interaction model with rejection option (SHIMR). A decision tree explains to a patient the diagnosis with a long rule (i.e., conjunction of many intervals), while SHIMR employs a weighted sum of short rules. Using proteomic data of 151 subjects in the Alzheimer's Disease Neuroimaging Initiative (ADNI) dataset, SHIMR is shown to be as accurate as other non-interpretability methods (Sensitivity, SP = 0.84 ± 0.13 and Area Under the Curve, AUC = 0.86 ± 0.09)."

"Only a small number of published studies performed a 'real' integration of omics and non-omics (DNO) data, mainly to predict cancer outcomes. Challenges in DNO data integration regard the nature and heterogeneity of non-omics data, the possibility of integrating large-scale non-omics data with high-throughput omics data, the relationship between DNO data (i.e., ascertainment bias), the presence of interactions, the fairness of the models, and the presence of subphenotypes. These challenges demand the development and application of new analysis strategies to integrate DNO data. In this contribution we discuss different attempts of DNO data integration in clinical and epidemiological studies. The integrative strategies used in the identified works adopted three modeling methods: Independent, conditional, and joint modeling."

"We set out to study if gene expression based predictors of chemotherapy resistance that are specific for breast cancer subtypes can improve upon the performance of generic predictors. For HER2+, ER+ breast cancer, subtype specific predictor based on clinical features outperformed the generic, non-specific predictor. This can be explained by the fact that the generic predictor included HER2 and ER status, features that are predictive over the whole set, but not within this subtype. In all other scenarios the generic predictors outperformed the subtype specific predictors or showed equal performance."

"The identification of robust lists of molecular biomarkers related to a disease is a fundamental step for early diagnosis and treatment. However, methodologies for the discovery of biomarkers using microarray data often provide results with limited overlap. These differences are imputable to 1) dataset size (few subjects with respect to the number of features); 2) heterogeneity of the disease; 3) heterogeneity of experimental protocols and computational pipelines employed in the analysis. In this paper, we focus on the first two issues and assess, both on simulated (through an in silico regulation network model) and real clinical datasets, the consistency of candidate biomarkers provided by a number of different methods. The simulated data allowed us to outline advantages and drawbacks of different methods across multiple studies of varying number of samples and to evaluate precision of feature selection on a benchmark with known biomarkers. Although comparable classification accuracy was reached by different methods, the use of external cross-validation loops is helpful in finding features with a higher degree of precision and stability."

"The Pan-Cancer Analysis Project aimed to identify the genomic changes in cancer types from the Cancer Genome Atlas (TCGA). The meaning of architectural features in clear cell renal cell carcinoma (ccRCC) by Fuhrman grade (FG) has not been investigated at clinic-pathologic or genetic levels in this set. Clinical data were also collected (gender, age, and stage). We used a Random Forest machine learning approach comparing low FG (1-2, 174 cases) vs. high FG (3-4, 139 cases) grade. The age, gender, and variants model performed with an AUC of 0.83. Here we report development, clinical validation, and diagnostic accuracy of a pre-operative molecular test (Affirma BRAF) to identify BRAF V600E mutations using mRNA expression in thyroid fine needle aspirate biopsies (FNABs). The resulting 238-gene linear support vector machine was compared to qPCR in the independent test set. Clinical sensitivity and specificity for malignancy were evaluated in a subset of test set samples (n=213) with expert-derived histopathology. We observed high positive (PPA, 90.9%) and negative (NPA, 99.0%) percent agreement with qPCR on the test set. Clinical sensitivity for malignancy was 43.8% (consistent with published prevalence of BRAF V600E in this neoplasm) and specificity was 100%, identical to qPCR on the same samples."

"In this study, machine learning methods (e.g., deep learning) were used to identify informative features from genome-scale omics data and to train classifiers for predicting the effectiveness of drugs in cancer cell lines. The methodology introduced here can accurately predict the efficacy of drugs, regardless of whether they are molecularly targeted or nontargeted chemotherapy drugs. This approach, on a per-drug basis, can identify sensitive cancer cells with an average sensitivity of 0.82 and specificity of 0.82, on a per-cell line basis, can identify effective drugs with an average sensitivity of 0.80 and specificity of 0.82."

"In this paper, we are interested in the question of how many and which genes should be selected for a disease class prediction. Our work consists of a Bayesian supervised statistical learning approach to refine gene signatures with a regularization which penalizes for the correlation between the variables selected. Our novel Bayesian approach includes a prior which penalizes highly correlated features in model selection and is able to extract key genes in the highly correlated context of microarray data. On real microarray datasets, we show that our approach can refine gene signatures to obtain either the same or better predictive performance than other existing methods with a smaller number of genes."

"The manuscript reviews several commonly used approaches to evaluate classification and clustering methods (e.g. cross-validation and bootstrap resubstitution). The identification of genomic biomarkers is a key step towards improving diagnostic tests and therapies. We present a reference-free method for this task that relies on a k-mer representation of genomes and a machine learning algorithm that produces intelligible models. The method was validated by generating models that predict the antibiotic resistance of *C. difficile*, *M. tuberculosis*, *P. aeruginosa*, and *S. pneumoniae* for 17 antibiotics. The obtained models are accurate, faithful to the biological pathways targeted by the antibiotics, and they provide insight into the process of resistance acquisition. The method is not limited to predicting antibiotic resistance in bacteria and is applicable to a variety of organisms and phenotypes."

"GEP NENs (gastroneuroendocrine pancreatic neuroendocrine neoplasms) were investigated by reverse-engineering intracellular signaling networks and identifying hub genes using degree (number of interactions) and betweenness (number of shortest paths) statistics. A random forest algorithm was used to assess hub gene expression in 130 blood samples (NENs: n = 63) and to differentiate healthy controls and GEP NENs. Gene-based classifiers detected NENs in independent sets with high sensitivity (85-98%), specificity (93-97%), PPV (95-96%) and NPV (87-98%). Additionally, multi-transcript logistic machine learning algorithms substantially outperform single-analyte assays such as CgA."

"The goal of the study was to identify severe asthma exacerbation children using phenotypic and SNP data. We concluded that the new classifier with the Newton-Raphson iterative processes and propensity scores have reliable performance with the increase in AUC values in all cases: (i) phenotypic data only; (ii) phenotypic data with the top ten significant SNPs; and (iii) phenotypic data with the top 160 to 302 significant SNPs."

"Current imaging modalities are inadequate in preoperatively predicting regional lymph node metastasis (RLNM) status in rectal cancer (RC). Here, we designed support vector machine (SVM) model to address this issue by integrating epithelial-mesenchymal transition (EMT)-related biomarkers along with clinicopathological variables. The sensitivity, specificity and overall accuracy of SVM in predicting RLNM were 68.2%, 81.1% and 72.3%, respectively. Importantly, multivariate logistic regression analysis showed that SVM model was indeed an independent predictor of RLNM status (odds ratio, 11.536; 95% confidence interval, 4.132-36.36; P<0.0001)."

"Drug response prediction is a critical step for personalized treatment of cancer patients and ultimately leads to precision medicine. In this paper, we proposed a method based on quantile regression forest and applied it to the CCLE dataset. Through the out-of-bag validation, our method achieved much higher prediction accuracy of drug response than other available tools."

"Urine proteome analysis (UPA) has already provided accurate discriminatory patterns of urinary peptides for renal disease, coronary artery disease, and asymptomatic LV diastolic dysfunction. UPA has now been used to characterize a discriminatory peptide biomarker pattern and establish a diagnostic classifier for heart failure patients with reduced ejection fraction (HFrEF) in the presence of chronic kidney disease (CKD). In total, 107 significant discriminatory peptides were identified and used to establish a support vector machine-based classifier that was successfully applied to a test set of 25 HFrEF patients and 33 controls, achieving 84% sensitivity and 91% specificity."

Author(s)	Title	Journal	Year	Country	Study Design	Key Findings	Conclusion
Fasbender, A and Waekens, E and Kyama, C and Bokor, A and Vodolazkaja, A and Verbeek, N and Van De Plas, R and Ojeda, F and Gevaert, O and Meuleman, C and Peeraer, K and Tomassetti, C and De Moor, B and D'Hooghe, T	Biomarkers in plasma or serum: Pitfalls in data processing	Journal of Clinical Oncology	2018	Belgium	Case-control study	254 plasma samples from women with (n=165) & without (n=89) endometriosis	accuracy (data were divided randomly (100 times) into training set (70%) and test set (30%))
Filmore, N and Ramos-Cejudo, J and Cheng, D and Turk, D and Shiels, A R and Chen, D and Elibers, D and Sung, F C and Johnson, B and Shannon, C and Pierce-Murray, K and Gaynor, K and Demedoneo, S C and Schiller, S and Ajjarapu, S and Hall, R and Ayvazian, S and Meng, F and Briphy, M T and Do, N	A Predictive model for survival in non-small cell lung cancer (NSCLC) based on genomic, phenotypic, clinical, and tumor sequencing data at the Department of Veterans Affairs (VA)	Journal of Clinical Oncology	2019	USA	Case-control study	356 VA patients newly diagnosed with NSCLC	Precision, recall, and area under the ROC curve (AUC) (5-fold CV)
55 Firoozbakhsh, F and Rezaei, A and D'Agostino, M and Porter, L and Ruveda, L and Nigam, A	An Integrative Approach for Identifying Network Biomarkers of Breast Cancer Subtypes Using Genomic, Interactomic, and Transcriptomic Data	J Comput Biol	2017	Canada	Tumour subtype categorization	"We have used the METABRIC data set (Curtis et al., 2012), which contains the copy number values and G levels of 2000 primary breast tumors with long-term clinical follow-up"	AUC, sensitivity, specificity (10-fold CV)
Fong, F and Bar, H Y and Shedd, K and Salya-Cork, K and Oullette, P and Campagne, F and Melnick, S E and Mahb, S and Shakhovich, R	Epigenetic profiling of primary CLL reveals novel DNA methylation-based clusters and novel mechanisms of leukemogenesis	Blood	2012	USA	Case-control study	"DNA methylation of over 240 patients with CLL"	AUC (10-fold CV)
57 Gal, O and Auslander, N and Fan, Y and Meerzaman, D	Predicting Complete Remission of Acute Myeloid Leukemia: Machine Learning Applied to Gene Expression	Cancer Informatics	2019	USA	Case-control study	473 bone marrow specimens from 473 patients	AUC (5-fold CV + test set)
58 Gamberger, D and Lavrac, N and Zelency, F and Tolar, J	Induction of comprehensible models for gene expression datasets by subgroup discovery methodology	J Biomed Inform	2004	Croatia	Case-control study	the approach was applied to multiple cancer microarray datasets with ~50 samples per group in total	sensitivity, specificity, precision (training/test set split)
59 Gao, H and Zheng, Z and Yue, Z and Liu, F and Zhou, L and Zhao, X	Evaluation of serum diagnosis of pancreatic cancer by using surface-enhanced laser desorption/ionization time-of-flight mass spectrometry	Int J Mol Med	2012	China	Case-control study	serum samples from 132 patients with PCA and 67 healthy controls	sensitivity, specificity (leave-one-out cross-validation)
Gerl, M J and Klose, C and Surma, M A and Fernandez, C and Melander, O and Mannisto, S and Borodulin, K and Havulinna, A S and Salonen, V and Ikonen, Y and Cammistrà, C V and Simons, K	Machine learning of human plasma lipoproteins for obesity estimation in a large population cohort	PLoS Biology	2019	China	Case-control study	Samples of the FHSR2012 underwent lipoprotein measurements (1,214 randomly selected individuals) of which 1,065 were used	R-squared of obesity indicator variables (5x repeated 10-fold CV)
61 Gong, J and Fox, N S and Huang, V and Boutros, P C	Prediction of early breast cancer patient survival using ensembles of hypoxia signatures	PLoS One	2018	Canada	Cases only (survival prediction)	1,546 early breast cancer patients	AUC, accuracy, sensitivity, specificity (10-fold cross-validation + test set)
62 Gram, K and Friedl, V and Houhain, K E and Stuart, J M	PLATYPUS: A Multiple-View Learning Predictive Framework for Cancer Drug Sensitivity Prediction	Pac Symp Biocomput	2019	USA	Cases only (drug sensitivity prediction)	At the time of download the Cancer Cell Line Encyclopedia (CCLE) contained genomic, phenotypic, clinical, and other annotation data for 1,037 cancer cell lines	cross-validation + external cohort validation
63 Grellemund, V and Prasad, P F and Querin, G and Debot, F and Le Chat, G and Prasad-Peyre, J F and Bede, P	Machine learning in amyotrophic lateral sclerosis: Achievements, pitfalls, and future directions	Frontiers in Neurosci	2019	France	review (not applicable)	"The first dataset we used was collected from Genomic of Drug Sensitivity in Cancer project (release 5.0, https://www.cancerrgenet.org/download), including 652 cancer cell lines, 135 drugs, and 70,676 known response values. The second dataset was collected from the CCLE (https://portals.broadinstitute.org/ccle), which contains 23 drugs and 493 cell lines with 10,870 known responses"	review
64 Guan, N N and Zhao, Y and Wang, C and Liu, Q and Chen, X and Piao, X	Molecular Therapy + Nucleic Acids	Journal of Clinical Oncology	2019	China	Cases only (drug response prediction)	"Serum samples were collected from all participants, including 587 CPO patients (CX2 + 120, CX22 + 104, CX23 + 110, CX24 + 102) and 100 healthy controls"	Pearson correlation coefficient (PCC), root-mean-square error (RMSE), PCCr, and RMSEr averaged over all drugs (10-fold CV)
65 Guo, S and Guo, D and Chen, L and Jiang, Q	A centroid-based gene selection method for microarray data classification	J Theor Biol	2016	China	Case-control study	multiple microarray datasets for different cancers with ~50 sample per group in total were used	accuracy + standard deviation (repeated 5-fold CV)
66 Guo, Y and Yu, H and Chen, D and Zhao, Y Y	Machine learning distilled metabolite biomarkers for early stage renal injury	Metabolomics	2020	China	Case-control study	"Serum samples were collected from all participants, including 587 CPO patients (CX2 + 120, CX22 + 104, CX23 + 110, CX24 + 102) and 100 healthy controls"	AUC (10-fold CV)

Author	Title	Journal	Year	Country	Study Type	Methodology	Findings	Classification
67 Han, H	Nonnegative principal component analysis for mass spectral serum profiles and biomarker discovery	BMC Bioinformatics	11	USA	2010	Case-control study	4 different cases/control cancer MS serum profile datasets were analyzed with >50 samples per group on average	accuracy (LOOCV and 100 trials of 50% holdout cross validations (HOCV))
68 Fedorovici, G and Walsh, P S and Sadov, P M and Huang, J and Kennedy, G C	Identifying of Hürthle cell cancers: solving a clinical challenge with genomic sequencing and a use of machine learning algorithms	Bmc Systems Biology	13	USA	2019	Case-control study	318 samples, including 119 Hürthle cell-negative and 199 Hürthle cell-positive samples	AUC, sensitivity, specificity (10-fold nested CV)
69 Heard, B J and Rosovick, J M and Fritzer, M J and Edl-Gabslaway, H and Wiley, J P and Krawetz, R J	A computational method to differentiate normal individuals, osteoarthritis and rheumatoid arthritis patients using serum biomarkers	J R Soc Interface	11	Canada	2014	Case-control study	normal individuals (normal, n = 100), patients with osteoarthritis (OA, n = 100), and rheumatoid arthritis (RA, n = 100)	accuracy, sensitivity, specificity (training, validation and test set split)
70 Hernández, B and Pennington, S R and Parnell, A C	Bayesian methods for proteomic biomarker development	EuPA Open Proteomics	9	USA	2015	review (not applicable)	review	review
71 Hiss, Z M and Gilles, D F	Identifying significant features in cancer methylation data using gene pathway segmentation	Cancer Informatics	15	UK	2016	Case-control study	normal individuals (normal, n = 100), patients with osteoarthritis (OA, n = 100), and rheumatoid arthritis (RA, n = 100)	AUC, accuracy (stratified 10-fold CV)
72 Ho, D S and Wand Schierding, W and Wake, M and Saffery, R and O'Sullivan, J	Machine learning SNP based prediction for precision medicine	Frontiers in Genetics	10	Australia	2019	review (not applicable)	review	review
73 Honda, K and Kashiwagi, Y and Umata, T and Okuyama, T and Kizugawa, T and Kikuchi, S and Endo, M and Tsuchida, A and Aoyagi, T and Ito, T and Morioka, F and Hirohata, S and Yamada, T	Possible detection of pancreatic cancer by plasma protein profiling	Cancer Res	65	Japan	2005	Case-control study	71 pancreatic cancer patients and 71 healthy controls	AUC, sensitivity, specificity (LOOCV, external test set)
74 and Choi, A M and Baron, R M and Thomas, N J and Wong, H R and Broadie, J R and Chonhchi, V M	Genomics analysis of gene expression and CRISPRi demonstrates a distinct ARDS signature	American Journal of Respiratory and Critical Care Medicine	197	USA	2018	training cohort (n = 318, 75%), validation cohort (n = 105, 25%)	training + test set	training + test set
75 and McDonald, J F	Machine learning predicts individual cancer responses to therapeutic drugs with high accuracy	Sci Rep	8	USA	2018	175 cancer patients	accuracy, sensitivity, specificity (LOOCV)	cross-validation
76 Ma, H and Hou, C C	Predicting Breast Cancer by Paper Spray Ion Mobility Spectrometry Mass Spectrometry and Machine Learning	Analytical Chemistry	92	USA	2020	breast core needle biopsies: 29+177 benign, 14+0 malignant	accuracy, sensitivity, specificity (cross-validation + external validation)	cross-validation + external cohort validation
77 Huang, Y H and Kuo, H C and Li, S C and Chiu, X Y and Liu, S F and Kuo, H C	HAMP promotes hypomethylation and increased promoter levels as biomarkers for Kawasaki disease	J Mol Cell Cardiol	117	England	2018	241 cases, including 128 KD patients, who were tested both prior to receiving intravenous immunoglobulin (IVIG) and at least 3 weeks after IVIG treatment, and 128 fibrin controls, who were observed in the Illumina HumanMethylation450 BeadChip study for their CpG markers. The remaining cases consisted of another 92 KD patients and 113 controls that were used for validation by pyrosequencing	AUC, sensitivity, specificity (5-fold CV, external test set)	cross-validation + external cohort validation
78 and Li, Y Z	Proximal diagnostic model for systemic lupus erythematosus using proteomic fingerprint technology	Shkhan Da	40	China	2009	430 patients (218 SLE, 122 non-SLE)	sensitivity, specificity (training/test set split)	training + test set

1	Ishii, H and Saitoh, M and Sakamoto, K and Sakamoto, K and Saigusa, D and Kasa, H and Ashizawa, K and Miyazawa, K and Taketsu, S and Masuyama, K and Yoshimura, K	Uptime-based rapid diagnosis with machine learning for detection of TGF- β signaling activated area in head and neck cancer	10-10	Japan	https://doi.org/10.1136/tj-2014-14-002-07	article	A total of 240 and 90 mass spectra were obtained from TGF- β stimulated and non-stimulated HNSCC cells, respectively	Case-control study	accuracy (LOOCV)	cross-validation	"We established a rapid diagnostic system based on the combination of probe electrospray ionisation-mass spectrometry (PESI-MS) and machine learning without the aid of immunohistological and biochemical procedures to identify tumour areas with heterogeneous TGF- β signalling status in head and neck squamous cell carcinoma (HNSCC). This discriminant algorithm achieved 98.79% accuracy in discrimination of TGF- β stimulated cells from untreated cells." "A major obstacle to reducing the mortality of colorectal cancer (CRC) is prompt detection and treatment, and a simple blood test is likely to have higher compliance than all of the current methods. The purpose of this report is to examine the utility of a mass spectrometry-based blood serum protein biomarker test for detection of CRC. A five-marker panel consisting of leucine-rich alpha-2-glycoprotein 1, epidermal growth factor receptor, inter-alpha-trypsin inhibitor heavy chain family member 4, hemopectin, and superoxide dismutase 1 was performed the best with 70% specificity at over 80% sensitivity (area under the curve = 0.86) in the validation set."				
2	Ivancic, M M and Megna, B W and Serchukov, Y and Craven, M and Reichelderfer, M and Pickhardt, P J and Susman, M R and Kennedy, G D	Noninvasive Detection of Colorectal Adenomas Using Serum Protein Biomarkers	J Surg Res	246	160-169	2020	United States	https://doi.org/10.1097/JRS.0000000000000718	article	"Blood was drawn from individuals (n = 153) before colonoscopy or from patients with nonmetastatic CRC (n = 50). AML dataset: "194 samples contain methylation data and we use the part of the data measured by NHI-USC HumanMethylation450 arrays. 173 samples contain mRNA data measured by Hi-C113 arrays." BRCA dataset: "This data set includes 993 samples with clinical data. Only very few samples in this data set are indicated as having metastasized (8 samples). Hence the data are analyzed according to "tumour size", "affected nearby lymph nodes", "stage", and "estrogen receptor"."	Case-control study	AUC, sensitivity, specificity (training / test set split)	training + test set	"Molecular measurements from cancer patients such as gene expression and DNA methylation can be influenced by several external factors. If a model does not take potential biases in the data into account, this can lead to problems when trying to predict the stage of a certain cancer type. This is especially true when these biases differ between the training and test set. We introduce a method that can estimate this bias on a per-feature level and incorporate calculated feature confidences into a weighted combination of classifiers with disjoint feature sets. Moreover, we show how to visualize the learned classifiers to display interesting associations with the target label."	
3	82 Jalali, Ad and Pfeifer, N	Interpretable per case weighted ensemble method for cancer associations	BMC Genomics	17	501-501	2016	Germany	https://doi.org/10.1186/s12864-016-0874-4	article	Case only (risk & severity stratification)	AUC (training / test set split)	training + test set	"Molecular measurements from cancer patients such as gene expression and DNA methylation can be influenced by several external factors. If a model does not take potential biases in the data into account, this can lead to problems when trying to predict the stage of a certain cancer type. This is especially true when these biases differ between the training and test set. We introduce a method that can estimate this bias on a per-feature level and incorporate calculated feature confidences into a weighted combination of classifiers with disjoint feature sets. Moreover, we show how to visualize the learned classifiers to display interesting associations with the target label."		
4	Jones, JJ and Wilk, B E and Benz, R W and Babbar, N and Borgagine, G and Burrell, T and Christie, E B and Cronin, C and Cole, P and Dillon, R and Kirs, J R and Kuo, A and Prestrel, J and Schroeder, S R and Sklar, H and Smith, W F and You, J and Hills, D W and Agus, D B and Blume, J E	A Plasma-Based Protein Marker Panel For Colorectal Cancer Detection Identified by Multiplex Targeted Mass Spectrometry	Clin Colorectal Cancer	15	2	186-194	2016	United States	https://doi.org/10.1016/j.ccr.2016.04.004	article	the present study used 274 individuals patient blood plasma samples, 137 with biopsy confirmed colorectal cancer and 137 age- and gender-matched controls.	Case-control study	AUC, sensitivity, specificity (cross-validation + external test)	validation	"Early colon cancer detection in patient populations ineligible for testing, such as the elderly or those with significant comorbidities, could have clinical benefit. A multiplex assay was developed for 317 candidate marker proteins, using 337 peptides monitored through 674 simultaneously measured MRM transitions in a 30-minute liquid chromatography-mass spectrometry analysis of immunodepleted blood plasma. To evaluate the combined candidate marker performance, the present study used 274 individual patient blood plasma samples, 137 with biopsy-confirmed colorectal cancer and 137 age- and gender-matched controls. Using one half of the data as a discovery set (69 disease cases and 69 control cases), the elastic net feature selection and Random Forest classifier assembly were used in cross-validation to identify a 15-transition classifier. The mean training receiver operating characteristic area under the curve was 0.82. After final classifier assembly using the entire discovery set, the 136-sample (69 disease cases and 68 control cases) validation set was evaluated. The validation area under the curve was 0.91."
5	Jurmeister, P and Bockmayr, M and Seegerer, F and Bockmayr, T and Treue, D and Montavon, G and Vollenbreck, C and Arnold, A and Teichgraber, D and Bressan, G and Schuller, U and von Laffert, M and Muller, K R and Capper, D and Klauschen, F	Machine learning analysis of DNA methylation profiles distinguishes primary lung squamous cell carcinoma from head and neck metastases	Science Translational Medicine	11	509	10-10	2019	Germany	https://doi.org/10.1126/scitranslmed.aau1113	article	408 patients with a history of primary HNSCC and a synchronous or metachronous squamous lung primary	Case-control study	AUC, accuracy (5-fold CV + external test)	cross-validation + external cohort validation	"Head and neck squamous cell carcinoma (HNSC) patients are at risk of suffering from both pulmonary metastases or a second squamous cell carcinoma of the lung (LUSC). Differentiating pulmonary metastases from primary lung cancers is of high clinical importance, but not possible in most cases with current diagnostics. To address this, we performed DNA methylation profiling of primary tumors and trained three different machine learning methods to distinguish metastatic HNSC from primary LUSC. We developed an artificial neural network that correctly classified 98.4% of the cases in a validation cohort of 279 patients with HNSC and LUSC as well as normal lung controls, outperforming support vector machines (95.7%) and random forests (87.8%)."
6	85 Karimpour-Fard, A and Epperson, L E and Hunter, L E	A survey of computational tools for downstream analysis of proteomic and other omics datasets	Human Genomics	9	11-11	2015	USA	https://doi.org/10.1186/s12916-015-0005-2	article	review (not applicable)	review	"In this paper, we review the well-known and ready-to-use tools for classification, clustering and validation, interpretation, and generation of biological information from experimental data. We suggest some rules of thumb for the reader on choosing the best suitable learning method for a particular dataset and conclude with pathway and functional analysis and then provide information about submitting final results to a repository."			
7	Khan, S R and Mohan, H and Liu, Y and Bachchu, B and Gehli, J and Al Rajjal, D and Manialayal, Y B and Cox, B J and Gunderson, E A and Wheeler, M B	The discovery of novel predictive biomarkers and early-stage pathobiology for the transition from gestational diabetes to type 2 diabetes	Diabetologia	6	2	687-703	2019	Canada	https://doi.org/10.1007/s00125-019-04802-z	article	55 incident cases matched to 85 non-case control participants	Case-control study	AUC, accuracy, sensitivity, specificity (45-fold cross-validation)	cross-validation	"Gestational diabetes mellitus (GDM) affects up to 20% of pregnancies, and almost half of the women affected progress to type 2 diabetes later in life, making GDM the most significant risk factor for the development of future type 2 diabetes. We used a well-characterized prospective cohort of women with a history of GDM pregnancy, all of whom were enrolled at 6-9 weeks postpartum (baseline), were confirmed not to have diabetes (n = 75) OGTT and tested annually for type 2 diabetes on an ongoing basis (2 years of follow-up). A large-scale targeted (epidemic study) was implemented to analyze >1100 lipid metabolites in baseline plasma samples [...] Machine learning optimization in a decision tree format revealed a seven lipid metabolite type 2 diabetes predictive signature with a discriminating power (AUC) of 0.92 (87% sensitivity, 93% specificity and 91% accuracy)."
8	Khusial, R D and Cioffi, C E and Calhoun, S A and Kasinakis, A M and Alzarak, A and Knight-Scott, J and Aletto, R and Castillo-Leon, E and Jones, D P and Pierpont, J and Caprio, S and Santoro, N and Kim, A and You, M B	Development of a Plasma Screening Panel for Pediatric Nontoxicology Fatty Liver Disease Using Metabolomics	Hepatology Communications	3	10	1311-1321	2019	United States	https://doi.org/10.1002/hep.4147	article	subjects with NAFLD (n = 222) and without NAFLD (n = 337)	Case-control study	AUC (training set: 2/3 of data, test set: 1/3 of data)	training + test set	"Nonalcoholic fatty liver disease (NAFLD) is the most common chronic liver disease in children, but diagnosis is challenging due to limited availability of noninvasive biomarkers. Machine learning applied to high-resolution metabolomics and clinical phenotype data offers a novel framework for developing a NAFLD screening panel in youth. Here, untargeted metabolomics by liquid chromatography-mass spectrometry was performed on plasma samples from a combined cross-sectional sample of children and adolescents ages 2-25 years old with NAFLD (n = 222) and without NAFLD (n = 337) [...]. The highest performing classification model was random forest, which had an area under the receiver operating characteristic curve (AUROC) of 0.94, sensitivity of 73%, and specificity of 97% for detecting NAFLD cases. A second classification model was developed using the metabolomic assessment of insulin resistance substituted for the WBV. Similarly, the highest performing classification model was random forest, which had an AUROC of 0.92, sensitivity of 73%, and specificity of 94%."
9	Kim, J and Du Ross, J C and Lee, J and Tomalin, L and Lowes, M A and Fitz, L and Bernstein, G and Valdez, A and Schwikowski, R and Krueger, J C and Sakara Santika, M	Precision medicine in psoriasis: Machine learning and proteomics join forces to develop a blood-based test for Enalapril in psoriasis patients	Experiment in Dermatology	25	49-50	2016	United States	https://doi.org/10.1111/exd.12020	meeting abstract	259 serum samples from a phase 3 study in adults with moderate-to-severe psoriasis	Case-control study	AUC, accuracy, learning based method for more accurate identification of prognostic biomarker genes and use them for prediction of cancer prognosis. The proposed method specifies the candidate prognostic gene module by graph learning using the generative adversarial networks (GANs) model, and scores genes using a PageRank algorithm. We applied the proposed method to multiple-omics data that included copy number, gene expression, DNA methylation, and somatic mutation data for five cancer types. The proposed method showed better prediction accuracy than did existing methods."	training + test set	"Despite various recent FDA-approved treatments 20%-30% of psoriatic patients fail to respond to biologics. [...] we applied machine-learning algorithms to proteomic data obtained using a proximity extension assay [...] to develop a blood-based test to predict response to tofacitinib or Enalapril in psoriasis patients. We elastic-net methods using pre-treatment data, was the best performer among the methods evaluated, with average AUC values (over 500 random 20/80 data splits) of 90.8% and 91.4% for accuracy of 88% and 87% for tofacitinib and Enalapril, respectively."	
10	Kim, M and Oh, I and Ahn, J	An improved method for prediction of cancer prognosis by network learning	Genes	9	10	1-11	2018	Switzerland	https://doi.org/10.3390/genes910178	article	"First, we downloaded gene mRNA data, CNV data, DNA methylation data, SNP data, and clinical data for PRAAD, BRCA, CRC, and S142 from The Cancer Genome Atlas (TCGA) (more than 50 samples per group for multiple datasets) [...] We apply the Meta-SVM methods to two real examples of idiopathic pulmonary fibrosis expression profiles (IPF; 221 samples in four studies of biopsy outcome (i.e., case and control)) and breast cancer expression profiles provided by The Cancer Genome Atlas (TCGA) including mRNA, copy number variation (CNV) and epigenetic DNA methylation	Case-control study	AUC, accuracy (10-fold CV)	cross-validation	"We propose a meta-analytic support vector machine (Meta-SVM) that can accommodate multiple omics data, making it possible to detect consensus genes associated with disease across studies. Experimental studies show that the Meta-SVM is superior to the existing meta-analysis method in detecting true signal genes. In real data applications, diverse omics data of breast cancer (TCGA) and mRNA expression data of lung disease (idiopathic pulmonary fibrosis; IPF) were applied. As a result, we identified gene sets consistently associated with the disease across studies."
11	90 Kim, S and Jhong, J H and Lee, J and Koo, Y J	Meta-analytic support vector machine Biostat for integrating multiple data	Mitochondria	10	1	2017	South Korea	https://doi.org/10.1186/s12916-017-0828-2	article	"We demonstrate application of Meta-SVM methods to three real omics examples of breast cancer expression profiles (1658 samples in seven studies), IPF expression profiles (IPF; 221 samples in six studies) and The Cancer Genome Atlas multi-cancer methylation profiles (TCGA, http://cancergenome.nih.gov/). "125 surgical lung biopsies from 86 patients. 58 samples were identified by the expert panel as usual interstitial pneumonia, 23 as non-specific interstitial pneumonia, 16 as hypersensitivity pneumonitis, four as sarcoidosis, four as organizing pneumonia, and 18 as respiratory bronchiolitis, two as atypical pneumonia, and 18 as infectious pneumonia."	Case-control study	sensitivity, specificity (cross-validation)	cross-validation	"The purpose of this article is to develop a meta-analytic top scoring pair (MetaTSP) framework that combines multiple transcriptomic studies and generates a robust prediction model applicable to independent test studies. We proposed two frameworks: (1) averaging TSP scores or combining P-values from individual studies, to select the top gene pairs for model construction. We applied the proposed methods in simulated data sets and three large-scale real applications in breast cancer, idiopathic pulmonary fibrosis and pan-cancer methylation. The result showed superior performance of cross-study validation accuracy and biomarker selection for the new meta-analytic framework."	
12	91 Kim, S and Lee, C W and Tsang, C G	MetaTSP: a meta-analytic top scoring pair method for robust cross-study validation of omics prediction analysis	Bioinformatics	32	13	1966-1974	2016	South Korea	https://doi.org/10.1093/bioinformatics/btw116	article	"125 surgical lung biopsies from 86 patients. 58 samples were identified by the expert panel as usual interstitial pneumonia, 23 as non-specific interstitial pneumonia, 16 as hypersensitivity pneumonitis, four as sarcoidosis, four as organizing pneumonia, and 18 as respiratory bronchiolitis, two as atypical pneumonia, and 18 as infectious pneumonia."	Case-control study	Youden index (5-fold cross-validation)	cross-validation	"Idiopathic pulmonary fibrosis is a progressive fibrotic lung disease that distorts pulmonary architecture, leading to hypoxia, respiratory failure, and death. Diagnosis is difficult because other interstitial lung diseases have similar radiological and histopathological characteristics. We aimed to develop a molecular test that distinguishes interstitial pneumonia from other interstitial lung diseases in surgical lung biopsy samples [...]. The microarray classifier was trained on 77 samples and was assessed in a test set of 48 samples, for which it had a specificity of 92% (95% CI = 100 and a sensitivity of 82% (64-95))."
13	Kim, S Y and Diggins, J and Pankratz, D and Huang, J and Pagan, M and Siny, N and Tom, E and Anderson, J and Choi, Y and Lynch, D A and Steele, M P and Flaherty, K R and Brown, K K and Farah, H and Bukstein, M I and Parbo, A and Selman, M and Wolters, P J and Nathan, S D and Colby, T V and Myers, J A and Katsenelson, L A and Pagan, C and Kennedy, C F	Classification of usual interstitial pneumonitis in patients with interstitial lung disease: assessment of a machine learning approach using high-dimensional transcriptional data	Respiratory Medicine	3	6	2016	Mexico	https://doi.org/10.1016/j.rmed.2016.04.011	article	"Idiopathic pulmonary fibrosis is a progressive fibrotic lung disease that distorts pulmonary architecture, leading to hypoxia, respiratory failure, and death. Diagnosis is difficult because other interstitial lung diseases have similar radiological and histopathological characteristics. We aimed to develop a molecular test that distinguishes usual interstitial pneumonia from other interstitial lung diseases in surgical lung biopsy samples [...]. The microarray classifier was trained on 77 samples and was assessed in a test set of 48 samples, for which it had a specificity of 92% (95% CI = 100 and a sensitivity of 82% (64-95))."	Case-control study	accuracy (5-fold cross-validation)	cross-validation	"Idiopathic pulmonary fibrosis is a progressive fibrotic lung disease that distorts pulmonary architecture, leading to hypoxia, respiratory failure, and death. Diagnosis is difficult because other interstitial lung diseases have similar radiological and histopathological characteristics. We aimed to develop a molecular test that distinguishes usual interstitial pneumonia from other interstitial lung diseases in surgical lung biopsy samples [...]. The microarray classifier was trained on 77 samples and was assessed in a test set of 48 samples, for which it had a specificity of 92% (95% CI = 100 and a sensitivity of 82% (64-95))."	

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46

Author	Title	Journal	Year	DOI	Abstract	Study Design	Outcome	Validation				
93 Kim, Y and Bismeljer, J and Zwart, W and Westels, L F A and Via, D J	Genomic data integration by WDN-PARAFAAC identifies interpretable factors for predicting drug sensitivity in vivo	Netherlands	2019	https://doi.org/10.1093/bioinformatics/btz116	1815 genes by 935 cell line	Cases only (drug sensitivity prediction)	AUC (10-fold CV)	cross-validation				
94 Kim, Y and Kim, D and Kim, S Y	Prediction of Acquired Taxane Resistance Using a Personalized Pathway-Based Machine Learning Method	Cancer Res Treat	2019 (South Korea)	https://doi.org/10.1158/1078-0432.CCR.19-137	more than 50 samples per group for most human cancer cell line datasets considered	Cases only (drug response prediction in vitro)	AUC (LOOCV)	cross-validation				
Kirilgoz, T and Kilic, S and Abali, Z Y and Yaman, A and Kavayuz, S B and Ehan, H and M. Turan, S and 95 Hallar, J and Sagiroglu, M S and Beneket, A and Gurun, T	Simplifying the interpretation of steroid metabolome data by a machine-learning approach	Research in Paediatrics	2019	https://doi.org/10.1155/2019/1689	500 healthy controls and 427 treatment-naïve children with a disorder of adrenal steroidogenesis	Case-control study	sensitivity, specificity (10-fold cross-validation)	cross-validation				
Kitazawa, H and Muramatsu, H and Murakami, N and Okuno, Y and Wakamatsu, M and Yoshida, T and Imaya, M and Yamamoto, A and Miwata, S and Narita, K and Hamada, M and Ichikawa, D and Taniguchi, R and Kawashima, N and Nishikawa, E and Naita, A and Nishio, N and Kojima, S and 96 Takahashi, Y	Genome-wide methylation analysis of DNA methyltransferase enzyme activity of methylation for stratification of patients with juvenile myelomonocytic leukemia	Blood	2019	https://doi.org/10.1182/blood-2019-07-877064	99 children (67 boys and 32 girls) with JMML	Case-control study	accuracy (training / test set split)	training + test set				
97 Kong, A and Azenkott, R	Statistical Applications in Genetics and Molecular Biology	16	1	13-30	2017	Germany	2017	https://doi.org/10.1159/00045-011	"A dataset of 238 MALDI colorectal mass spectra and two datasets of 216 and 253 SELDI ovarian mass spectra respectively were used to test our approach."	Case-control study	accuracy (LOOCV)	cross-validation
98 Krawiec, J and Lukasz, T	The feature selection bias problem in relation to high-dimensional gene data	Artif Intell Med	66	63-71	2016	Poland	2016	https://doi.org/10.1016/j.artmed.2016.11.005	seven microarray datasets with ~50 samples/group for multiple datasets were used	Case-control study	accuracy (double LOOCV)	cross-validation
Krittanawong, C and Bombak, A S and Baber, U and Bangalor, S and Meseiri, F H and Wilson Tang, 99 W H	Future Direction for Using Artificial Intelligence to Predict and Manage Hypertension	Curr Hypertens Rep	20	9	75-75	2018	Poland	2018	https://doi.org/10.1007/s12974-018-0878-x	review (not applicable)	review	
Kuo, C H S and Pavlidis, S and Loza, M and Barbaud, F and Rowe, S and A. Pandis, I and Rossio, C and 100 Wilson, S and Djukanovic, R and Stern, P and Chung, K F and Adcock, I M and Guo, Y	Asthma is a heterogeneous disease underlined by different inflammatory programs. A three-gene Th2 signature was associated with airway hyper-responsiveness, allergy markers and inhaled corticosteroid response (Woodruff et al. AIRCM 2009;180:388). We phenotyped asthma using a semi-supervised machine-learning approach to analyze gene expression profiles. Training of a gene model for the cluster using nearest shrunken centroids yielded an 8 signature set with 82% cross-validation accuracy.	American Journal of Respiratory and Critical Care Medicine	191		2015			https://doi.org/10.1164/rccm.2015.04.0878	"Subjects with moderate-to-severe asthma recruited in the U-BOPRED study underwent fiberoptic bronchoscopy for bronchial biopsy (91) and brush (105) samples"	Case-control study	accuracy (cross-validation)	cross-validation
101 Kurta, M B	BMC Bioinformatics	15	8-8	2014	Poland	2014	https://doi.org/10.1186/1471-2106-15-8	4 microarray datasets were used, one contained > 50 samples per group	Case-control study	error-rate (training / test set split)	training + test set	
Kuwabara, H and Washbuchi, A and Sova, R and Enomoto, M and Ishizaki, T and Tsuchida, A and 102 Nagakawa, Y and Katsumata, K and Sugimoto, M	Salivary metabolomics for colorectal cancer detection	Annals of Oncology	30	v46-v46	2019			https://doi.org/10.1093/annonc/mdz066	"231 subjects with CRC, 99 subjects with polyps, and 2272 subjects with healthy controls"	Case-control study	AUC (training / test set split)	training + test set
Lacroix Triki, M and Kempowsky-Hamon, T and Valle, C and Hedjazi, L and Lamarque, S and Trouilh, L and Puydieu, S and Mahidi, L and Daler, F and Fillion, J and Favre, G and Le Lam, N V and 103 Berre-Anton, V	Fuzzy logic selection as a new reliable tool to identify genes in breast cancer - the INNODIAG Study	Laboratory Investigation	93	S1A-S1A	2013			https://doi.org/10.1093/lin/lbn231	7 breast cancer microarray datasets + 151 consecutive invasive breast carcinomas	Case-control study	sensitivity, specificity, error rate (training / test set split)	training + test set
La, A and Panos, R and Marjanovic, M and Walker, M and Fuentes, E and Kapp, D S and Henner, W D and Buturovic, L and Miller, M H	A gene expression profile test that distinguishes ovarian from endometrial cancers	Journal of Clinical Oncology	30	15	2012			https://doi.org/10.1200/JCO.2011.20311	75 metastatic, poorly differentiated or undifferentiated primary FPE tumor specimens.	Differential diagnosis prediction	AUC (training / test set split)	training + test set

Author(s)	Year	Country	Journal	Volume	Issue	Pages	DOI	Abstract	Study Design	Outcome	Validation	Notes
Lawton, K A and Brown, M V and Alexander, D L and Wu, J J and Lawson, R and Jaffe, M and Milburn, M V and Ryals, J A and Bowser, R and Kudoczka, M E and Berry, J D	2014	England	Frontiers in Oncology	5	362-370	2014	https://doi.org/10.3389/fonc.2014.00831	172 patients recently diagnosed with ALS, 50 healthy controls, and 73 neurological disease mimics. The SLE compendium contained 15,497 gene expression measurements with observations from healthy control (n=140) samples, treatment-naïve SLE (n=1,290) samples, and SLE samples exposed to various treatments (n=126)	Case-control study	AUC, sensitivity, specificity (training / test set split)	training + test set	
106 Le, T T and Blackwood, N O and Taroni, J H and Fu, W and Breitenstein, M K	2018	USA	AMIA Annu Symp Proc	2018	1358-1367	2018	https://www.ncbi.nlm.nih.gov/pmc/articles/PMC6171426/	five microarray datasets were used, including datasets with > 50 samples per group	Case-control + treatment response	balanced accuracy (cross-validation + 20% hold-out test set)	cross-validation + test set	
107 Leclercq, M and Vittrant, B and Martin-Magniette, M L and Scott Boyer, M P and Pernin, D and Bergeron, A and Fradet, Y and Drost, A	2019	Canada	Frontiers in Genetics	10		2019	https://doi.org/10.3389/fgen.2019.00406	Large-scale automatic feature selection for biomarker discovery in high-dimensional omics data	Case-control study	accuracy (ACC), balanced error rate (BER), Matthews' correlation coefficient (MCC), area under the curve (AUC), sensitivity, specificity, Root Mean Squared Error (RMSE), Correlation Coefficient (CC) (10-fold CV)	cross-validation	
108 Lee, S S and Attwood, K and Roder, H and Asmellah, S and Meyer, K and Kakolyris, S and Oliveira, C and Roder, J and Grigoriou, J and Chelis, L and Lye, R and Mahalingam, D	2019		Cancer Research	79	13	2019	https://doi.org/10.1158/1078-0432.CCR.18-2630	An independent validation of a screening test for plasma spectrometry for detection of hepatocellular carcinoma	Case-control study	AUC (training and validation cohort)	external cohort validation	
109 Lin, X and Afari, B and Marchionni, L and Cope, L and Parmigiani, G and Naiman, D and Geman, D	2009	USA	BMC Bioinformatics	10	256-256	2009	https://doi.org/10.1186/1471-2106-10-256	The ordering of expression among a few genes can provide simple biomarkers and signal BRCA1 mutations: a cross-study validation, including a dataset with > 50 samples per group	Case-control study	accuracy, sensitivity, specificity (LOOCV, cross-study validation)	cross-validation	
110 Liu, F and Xing, L and Zhang, X and Zhang, X	2019	China	Genes (Basel)	10	6	2019	https://doi.org/10.3390/genes10090144	A Four-Pseudogene Classifier Identified by Machine Learning Serves as a Novel Prognostic Marker for Survival of Osteosarcoma	Cases only (survival prediction)	AUC (10-fold CV)	cross-validation	
111 Liu, L and Liu, Y and Liu, C and Zhang, Z and Du, Y and Zhao, H	2016	China	Mol Med Rep	14	4	2016	https://doi.org/10.3892/mmr.2016.6	Analysis of gene expression profile identifies potential biomarkers for atherosclerosis	Case-control study	AUC (5-fold CV)	cross-validation	
112 Liu, M and Cand Jambhidi, A and Venn, O and Fields, A P and Maher, M C and Camm, G and Amin, H and Gross, S and Bredino, J and Miller, M and Schellenberger, J and Kurtzman, K N and Fung, E T and Maddala, T and Ovard, G and Klein, E A and Spiegel, D R and Hartman, A N and Aravani, A and Seiden, M	2019		Journal of Clinical Oncology	37		2019	https://doi.org/10.1200/JCO.2019.37.15.6	Genome-wide cell-free DNA (cfDNA) methylation signatures and effect on tissue of origin (TOO) performance	Tissue-of-origin prediction	accuracy (training / test set split)	training + test set	
113 Liu, W T and Wang, Y and Zhang, J and Ye, F and Huang, X H and Li, B and He, Q Y	2018	China	Cancer Lett	425	43-53	2018	https://doi.org/10.1016/j.canlet.2018.03.043	A novel strategy of integrated microarray analysis identifies CENPA, CENK and CCK2B as differentially expressed biomarkers in lung adenocarcinoma	Case-control study	accuracy (LOOCV, external test validation)	cross-validation + external cohort validation	
114 Wittenberg, G and Ye, J	2016	Belgium	BMC Genomics	17	669-669	2016	https://doi.org/10.1186/s12854-016-0253-2	Metabonomic biosignature differentiates melancholic depressive patients from healthy controls	Case-control study	accuracy, sensitivity, specificity (10-fold CV)	cross-validation	
115 Long, N P and Jung, K H and Yoon, S J and Anh, N H and Nghi, T D and Kang, Y P and Van, H H and Min, S W	2017	Vietnam	Oncotarget	8	65	109456	2017	2022 cancer, 115 cervical intraepithelial neoplasia (CIN), and 105 normal samples	Case-control study	accuracy, sensitivity, specificity (10-fold CV, external test set)	cross-validation + external cohort validation	
116 Long, N P and Nghi, T D and Kang, Y P and Anh, N H and Kim, H M and Park, S K and Kwon, S W	2020	USA	Metabolites	10	2	2020	https://doi.org/10.3390/met10020061	Toward a standardized strategy of clinical metabonomics for the advancement of precision medicine	review	review	review	
117 Long, N P and Park, S and Anh, N H and Nghi, T D and Yoon, S J and Park, J H and Lim, J and Kwon, S W	2019	Vietnam	Int J Mol Sci	20	2	2019	https://doi.org/10.3390/ijms20020090	High Throughput Omics and Statistical Observations of Precision Biomarker Discovery and Validation of Novel Diagnostic Signatures in Colorectal Cancer	Case-control study	AUC, sensitivity, specificity (5-times repeated 10-fold CV, test set)	cross-validation + test set	
118 Kwon, S W	2018	South Korea	Metabolites	8	109-109	2018	https://doi.org/10.3390/met8040109	Systematic assessment of cervical cancer initiation and progression uncovers genetic panels for deep learning-based early diagnosis and proposes novel diagnostic and prognostic biomarkers	review	review	review	
119 Long, N P and Tran, S J and Anh, N H and Nghi, T D and Lim, D K and Hong, Y J and Hong, S S and Kwon, S W	2020	South Korea	Metabolites	10	8	2020	https://doi.org/10.3390/met10080460	A systematic review on metabonomics-based diagnostic biomarker discovery and validation in pancreatic cancer	review	review	review	

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46

Author	Title	Journal	Year	Country	Study Type	Key Findings	Classification
129 Matlock, K and De Niz, C and Rahman, R and Ghosh, S and Paj, R	Investigation of model stacking for drug sensitivity prediction	BMC Bioinformatics	19	USA	2018	"We segregate 50 training samples into our vertical and horizontal groups, build individual predictive model RF with 50 trees, build the stacking model using a set of 150 samples, and obtain the performance MSEs of candidate models on a set of 50 testing samples. We then add 2 training samples and retrain the MSEs. We repeat this process until the training set has a total of 150 samples. The entire process is replicated 100 times with randomly selected training, testing, and validation samples in every iteration."	Cases only (drug sensitivity prediction)
McCarthy, C and Shrestha, S and Ibrahim, N E and Van Kimmede, R and Gaggin, H K and Mukai, R and Maqaret, C A and Barnes, G and Rhyne, R and Garate, J M and Januzzi, J L	Performance of a clinical/proteomic panel to predict obstructive peripheral artery disease in patients with and without diabetes mellitus	European Heart Journal	39	Netherlands	2018	"354 patients undergoing peripheral and/or coronary angiography, performance of this diagnostic panel was assessed in patients with (N=94) and without DM (N=260) using Monte Carlo cross-validation"	Case-control study
McDermott, J E and Wang, J and Mitchell, H and Webb-Robertson, B and Hafem, R and Ramey, J and Roland, K D	Challenges in biomarker discovery: Combining expert insights with statistical analysis of complex omics data	Expert Opinion on Medical Diagnostics	7	USA	2013	"While some small children display early-life wheeze, fewer have diagnosable asthma. We aimed to identify a metabolomic signature of early-life asthma, which could be used as a diagnostic test or which would provide insight into the biochemical basis of asthma in young children."	review (not applicable)
132 McGeachie, M and Kelly, K S and Litonjaa, A A and Weiss, S T and Lasky-Su, J A	Network of year-3 metabolites indicative of early-life asthma	Journal of Allergy and Clinical Medicine	197	USA	2018	cohort of 411 three-year olds at high risk for asthma	Case-control study
Meikle, P and Tsorotes, D and Barlow, C and Weir, J and Macintosh, G and Barber, M and Gousey, B and Stern, L and Kowalyk, A and Havel, J and White, A and Durr, A and Duffy, S and Kingwell, B	Plasma lipidomic analysis of stable and unstable coronary artery disease	Atherosclerosis Supplements	11	Australia	2010	"203 participants (control, n = 80, stable CAD, n = 61, unstable CAD, n = 62). We used the data from the Genomics of Drug Sensitivity in Cancer project [3], which contains 439 cancer cell lines, each of them characterised by a set of genomic features (details in the next section). The characterisation is not complete for every cell line, and therefore we filtered out cell lines with more than 15 missing genomic features, which reduced the set of selected cell lines from 639 to 608. The dataset contains 131 drugs."	Case-control study
Menden, M P and Iorio, F and Garnett, M and McDermott, U and Benez, C and Ballesera, P J and Saez-Rodriguez, J	Machine Learning Prediction of Cancer Cell Sensitivity to Drug Based on Genomic and Chemical Properties	PLoS One	8	United Kingdom	2013	"Predicting the response of a specific cancer to a therapy is a major goal in modern oncology that should ultimately lead to a personalised treatment. We developed machine learning models to predict the response of cancer cell lines to drug treatment, quantified through IC50 values, based on both the genomic features of the cell lines and the chemical properties of the considered drugs. Models predicted IC50 values in a 4-fold cross-validation and an independent blind test with coefficient of determination R2 of 0.72 and 0.64 respectively."	Cases only (drug sensitivity prediction)
135 Midonkawa, Y and Tsuji, S and Takayama, T and Aburatani, H	Genomic approach towards personalized anticancer drug therapy	Pharmacogenomics	13	Japan	2012	review (not applicable)	review
Mohadesany, P and Yousefi, S and Angadi, M and Gutman, D A and Barnholtz-Sloan, J S and Velazquez Vega, J E and Brat, D J and Cooper, J A D	Predicting cancer outcomes from histology and genomics using convolutional networks	Proc Natl Acad Sci U S A	115	USA	2018	"15 accuracy measurements, including Harrell's C-index for measuring concordance between predicted risks and actual survival (Monte Carlo cross-validation)"	Cases only (prognosis study)
137 Modlin, J and Kidd, M and Drezdov, I and Bodek, L and Malczewska, A and Matar, S	Automated finger print blood genomic diagnosis of neuroendocrine tumors	Neuroendocrinology	108	USA	2019	whole blood samples from 46,000 NETs and controls	Case-control study
138 Mohammed, A and Elgert, G and Adamczak, J and Helliker, T	Identification of potential tissue-specific cancer biomarkers and development of cancer versus normal genomic classifiers	Oncotarget	8	USA	2017	"A total of 2,175 tissue samples, both normal and cancerous, were collected from nine distinct tissues: blood [956], breast [171], colon [105], gastric [333], head and neck [82], lung [542], prostate [66], thyroid [234], and tongue [87]"	Case-control study
139 F D	Patient-specific cancer genes contribute to recurrently perturbed pathways and establish therapeutic vulnerabilities in esophageal adenocarcinoma	Nat Commun	10	Italy	2019	"261 esophageal adenocarcinomas (EACs) + 107 additional EACs for validation"	Cases only (survival prediction)
Murugesan, K and Ingle, M and Schrock, A B and Ngo, N and Frampton, G M and Alexander, B M and Miller, A A and Bekai-Sabb, T and Altabek, L A and Rovee, J S and Ai, J M	Comprehensive genomic profiling (CGP) of mixed hepatocellular and cholangiocarcinomas (mHCC-CCA)	Annals of Oncology	30	USA	2019	"1369 mHCC, 3865 CCA and 44 mHCC-CCA (> 80 samples for the main conditions)"	Case-control study
Nakamura, M and Bax, H J and Scotto, D and Souiri, E A and Sallie, S and Harris, R J and Hammar, N and Wainwright, A and Ghosh, S and Maitra, A and Spicer, J F and Van Heemrijck, M and Josephs, D H and Lacy, K E and Tsika, S and Karagannis, S N	Immune mediator expression signatures are associated with improved outcome in ovarian carcinoma	Oncotarget	8	Sweden	2019	1,656 ovarian carcinoma patient tumors	Cases only (survival prediction)

"A significant problem in precision medicine is the prediction of drug sensitivity for individual cancer cell lines. [...] We explore the predictive performance of model stacking and the effect of stacking on the predictive bias and squared error. In addition we discuss the analytical underpinnings supporting the advantages of stacking in reducing squared error and inherent bias of random forests in prediction of outliers. The performance of individual and stacked models are compared. We note that stacking models built on two heterogeneous datasets provide superior performance to stacking different models built on the same dataset. It is also noted that stacking provides a noticeable reduction in the bias of our predictors when the dominant eigenvalue of the principal axes of variation in the residuals is significantly higher than the remaining eigenvalues."

"Peripheral artery disease (PAD) is a global health problem associated with significant morbidity and mortality. Patients with diabetes mellitus (DM) are at substantial risk of developing PAD, however its diagnosis is often delayed until advanced stages when complications arise. Using proteomics and machine learning, a sub-set of artificial intelligence, we recently described a biomarker clinical/proteomic panel to predict prevalent obstructive PAD (HAART PAD) in patients undergoing diagnostic peripheral angiography and/or coronary angiography. In this study, we sought to compare the accuracy of this clinical/proteomic panel for the diagnosis of PAD in patients with and without DM. [...] In patients with DM, the HAART PAD panel had excellent performance or prediction of peripheral stenosis >50%. The model had an area under the receiver operating characteristic curve of 0.85 for obstructive PAD; higher scores were associated with greater severity of angiographic stenosis."

"In this review, we will present examples of current practices for biomarker discovery from complex omics datasets and the challenges that have been encountered in deriving valid and useful signatures of disease. We will then present a high-level review of data-driven (statistical) and knowledge-based methods applied to biomarker discovery, highlighting some current efforts to combine the two distinct approaches. Effective, reproducible and objective tools for combining data-driven and knowledge-based approaches to identify predictive signatures of disease are key to future success in the biomarker field. We will describe our recommendations for possible approaches to this problem including metrics for the evaluation of biomarkers."

"While some small children display early-life wheeze, fewer have diagnosable asthma. We aimed to identify a metabolomic signature of early-life asthma, which could be used as a diagnostic test or which would provide insight into the biochemical basis of asthma in young children. [...] The BN methodology identified a Bayesian Network profile of Year-3 asthma with 21 metabolites achieving an Area Under the Receiver Operator Characteristic Curve (AUC) of 86.5%."

"Currently there is no means to measure instability in coronary artery disease (CAD). Modified ceramide and modified phosphatidylcholine species were shown to distinguish stable and unstable CAD. These newly identified biomarkers were measured together with known plasma lipids, including sphingolipids, sphingomyelins, and phospholipids, to establish a plasma lipid profile using electrospray ionization tandem mass spectrometry. [...] Multivariate analysis using a statistical machine learning approach combined with recursive feature elimination and multiple cross-validation iterations was applied for the creation of prediction models. Comparison of models with varying number of features showed that models with only 8 lipids were sufficient to provide a similar discrimination between stable and unstable cohorts (AUC = 0.75) while 16 lipids were sufficient to discriminate control from CAD patients (AUC = 0.94)."

"Predicting the response of a specific cancer to a therapy is a major goal in modern oncology that should ultimately lead to a personalised treatment. We developed machine learning models to predict the response of cancer cell lines to drug treatment, quantified through IC50 values, based on both the genomic features of the cell lines and the chemical properties of the considered drugs. Models predicted IC50 values in a 4-fold cross-validation and an independent blind test with coefficient of determination R2 of 0.72 and 0.64 respectively."

"Here, we review recent advances in the development of classification algorithms using microarray technology for prediction of anticancer sensitivity, discuss the availability of public methods for prediction models, and present data regarding the identification of potential responders to FOLFOX therapy using random forests algorithms."

"We developed a computational approach based on deep learning to predict the overall survival of patients diagnosed with brain tumors from microscopic images of tissue biopsies and genomic biomarkers. This method uses adaptive feedback to simultaneously learn the visual patterns and molecular biomarkers associated with patient outcomes. Our approach surpasses the prognostic accuracy of human experts using the current clinical standard for classifying brain tumors and presents an innovative approach for accurate, objective, and integrated prediction of patient outcomes."

"NET Transcripts from GEP NETs were defined (2005-2008) and in 2010, RNA-seq classifiers based on machine learning developed. From 2010-2015, we refined using 41 53 genes co-expressed in blood and tumor tissue (N=48) and established a liquid biopsy (NETest). Diagnostic metrics were >90% sensitivity and specificity. Clinical evaluation (n=8,000 patients) demonstrated the multigene assay to be an effective prognostic biomarker."

"Machine learning techniques for cancer prediction and biomarker discovery can hasten cancer detection and significantly improve prognosis. Recent "OMICS" studies which include a variety of cancer and normal tissue samples along with machine learning approaches have the potential to further accelerate such discovery. To demonstrate this potential, 2,175 gene expression samples from nine tissue types were obtained to identify gene sets whose expression is characteristic of each cancer class. Using random forests classification and ten-fold cross-validation, we developed nine single-tissue classifiers, two multi-tissue cancer-versus-normal classifiers, and one multi-tissue normal classifier. Given a sample of a specified tissue type, the single-tissue models classified samples as cancer or normal with a testing accuracy between 85.29% and 92.6%. Given a sample of non-specific tissue type, the multi-tissue bi-class model classified the sample as cancer versus normal with a testing accuracy of 97.89%. Given a sample of non-specific tissue type, the multi-tissue multi-class model classified the sample as cancer versus normal and as a specific tissue type with testing accuracy of 97.43%. Given a normal sample of any of the nine tissue types, the multi-tissue normal model classified the sample as a particular tissue type with a testing accuracy of 97.55%."

"The identification of cancer-promoting genetic alterations is challenging particularly in highly unstable and heterogeneous cancers, such as esophageal adenocarcinoma (EAC). Here we describe a machine-learning algorithm to identify cancer genes in individual patients considering all types of damaging alterations simultaneously. [...] Experimentally mimicking the alterations of predicted helper genes in cancer and precursor cells validates their contribution to disease progression."

"mHCC-CCA is a rare primary liver carcinoma with histologic features of both hepatocellular carcinomas (HCC) and liver cholangiocarcinomas (CCA). In order to elucidate their shared and distinctive biology, we used Comprehensive genomic profiling (CGP) to compare the genomic alterations (GAs) of mHCC-CCA with those of HCC and CCA. [...] Utilizing the significant differences in GAs, biomarkers, and demographics between HCC and CCA (Fisher exact test, FDR < 0.05, only GAs with freq > 10% shown. Table), we built a random forest based machine learning model to rank a mHCC-CCA specimen in the CCA-HCC spectrum (Out of Bag error rate = 12.6%, AUC = 0.94). Biomarkers exclusively associated with one disease type included D1H1 (15% in CCA vs 1% in HCC, p=1e-57), TERT (4% in CCA vs 47% in HCC, p=2e-273) and Hepatitis B virus (HBV, 0.8% in CCA vs 8.5% in HCC, p=8e-42)."

"Immune and inflammatory cascades may play multiple roles in ovarian cancer. We aimed to identify relationships between expression of immune and inflammatory mediators and patient outcomes. We integrated differential gene expression of 44 markers and marker combinations (n = 1,978) in 1,656 ovarian carcinoma patient tumors, alongside matched 5-year overall survival (OS) data in silico. Expression of the 44 markers could discriminate between malignant and non-malignant tissues with at least 96% accuracy."

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46

142	Nakariyuki, S	A hybrid gene selection algorithm based on interaction information for microarray-based cancer classification	PLoS One	14	2	e0212333	2019	Thailand	https://doi.org/10.1371/journal.pone.0212333	article	ten per group data sets with > 50 samples microarray data multiple classes	Case-control study	accuracy, precision, recall, F-score (nested cross-validation)	cross-validation
143	Naorem, D and Muthaiyan, M and Venkatesan, A	Integrated network analysis and machine learning approach for the identification of key genes of triple-negative breast cancer	Journal of Cellular Biochemistry	120	4	6154-6167	2019	India	https://doi.org/10.1002/jcb.23904	article	Six microarray data sets consisting of 463 non-TNBC and 405 TNBC samples	Case-control study	AUC (training / test set split)	training + test set
144	Nazha, A and Komrokji, R S and Barnard, J and Ali-Ka, K and Padron, E and Madanat, Y F and Kumzomov, T and Akubaha, N and Steensma, D F and DeZern, A E and Roboz, G J and Garcia-Manero, G and Li, A F and Maciejewski, J P and Sekeres, M A	A Personalized Prediction Model to Risk Stratify Patients with Myelodysplastic Syndromes (MDS)	Blood	130			2017		https://doi.org/10.1182/blood-2017-09-240644	article	"Of 2302 pts, 1471 were included in the training cohort and 831 in the validation cohort"	Cases only (survival prediction)	C-index (training and validation cohort)	external cohort validation
145	Nazha, A and Sekeres, M A and Bejar, R and Komrokji, R S and Barnard, J and Ali-Ka, K and Prychodzen, B P and Hirsch, C M and Steensma, D P and DeZern, A E and Roboz, G J and Garcia-Manero, G and Ebert, B L and Maciejewski, J P	Genomic Biomarkers to Predict Response to Hypomethylating Agents in Patients with Myelodysplastic Syndromes (MDS)	Blood	130			2017		https://doi.org/10.1182/blood-2017-09-240644	article	433 pts with MDS (per 2008 WHO criteria) who received HMA (230 at our institution [training cohort], and 203 at other major academic institutions [validation cohort])	Cases only (treatment response prediction)	accuracy, sensitivity, specificity (training and validation cohort)	external cohort validation
146	Nozak, C and Calouso, A and Özgür, C J and Nyström, F H and Alam, M and Fedræk, T R and Kumzomov, T and Carroer Rog, J and Lepert, J and Heberich, P D and Cordero, A C and Lind, L and Ingelsson, E and Fall, T and Årnäs, J	Major proteomics for prediction of major cardiovascular events in type 2 diabetes	Diabetologia	61		S65-S65	2018	Brazil	https://doi.org/10.1007/s00125-018-0830-0	meeting abstract	1,211 adults with type 2 diabetes	Case-control study	accuracy with 95%-confidence intervals (training / test set split)	training + test set
147	O'Reilly, P and Orntutz, C and Gernon, G and O'Connell, E and Seahighe, C and Boyce, S and Serrano, L and Szegedi, E	Co-exing gene networks predict TRAIL responsiveness of tumour cells with high accuracy	BMC Genomics	15		1144-1144	2014	Ireland	https://doi.org/10.1186/s12854-014-0244-5	article	Gene expression microarray data for 109 tumor cell lines with known sensitivity to the death ligand cytokine tumor necrosis factor-related apoptosis-inducing ligand (TRAIL)	Case-control study	AUC, sensitivity, specificity (training and validation cohort)	external cohort validation
148	Oh, J H and Lotan, Y and Gurnani, P and Rosenblatt, K P and Gao, J	Prostate cancer biomarker discovery using high performance mass spectral serum profiling	Comput Methods Programs Biomed	96	1	33-41	2009	USA	https://doi.org/10.1016/j.cmpb.2008.04.004	article	Serum samples from 179 prostate cancer patients and 74 benign patients	Case-control study	accuracy, sensitivity, specificity, NPV, PPV (20 times 10-fold CV)	cross-validation
149	Okter, S and Pahlkaki, T and Aittokallio, T	Genetic variants and their interactions in disease risk prediction - Machine learning and network perspectives	BioData Mining	6	1		2013	Finland	https://doi.org/10.1186/1745-6216-6-3	article	review (not applicable) "The first cohort (EGAD00001001443, hereafter study cohort) contains RNAseq data and from CLL-purified cells of 136 individuals along with clinical data. The cohort was composed of 168 CLL, 22 monoclonal B cell lymphocytosis (MBL), and five small lymphocytic lymphoma (SLL) samples. There were 132 IGHV mutated cases and 64 IGHV unmutated cases in 119 males and 77 females. By staging at diagnosis, there were 22 MBL cases, 153 Binet Stage A cases, 18 Binet Stage B cases, and 8 Binet C stage cases. The second cohort (EGAD0001000258, hereafter validation cohort) consists of RNAseq data of CLL-purified cells from 98 individuals, of which 79 (81.7%) males and 24 (24.4%) females have publicly available phenotypic information. In this cohort there were 72 CLL, 4 SLL, and 3 MBL samples. 45 of the patients had mutated IGHV and 34 had unmutated IGHV. By staging at diagnosis, there were 3 MBL, 72 Binet Stage A, 3 Binet Stage B, and 1 Binet Stage C cases."	review		
150	Orpeiro, A M and Rodriguez, B A and Vence, N A and López Á, B and Arias, J A D and Varela, N D and Pérez, M S G and Encinas, M F P and López, L B	Time to treatment prediction in chronic lymphocytic leukemia based on new transcriptional patterns	Frontiers in Oncology	9			2019	Spain	https://doi.org/10.3389/fonc.2019.01607	article	Genetic variants and their interactions in disease risk prediction - Machine learning and network perspectives	Case-control study	accuracy, precision, recall, ROC (training and validation cohort)	external cohort validation
151	Ornel, J and Valjejo, E E and Estrada, K and Pena, J G and Alzheimer's Dis Neuroimaging, Initia	Benchmarking machine learning models for late-onset Alzheimer's disease prediction from genomic data	Bmc Bioinformatics	20	1	17-17	2019	Mexico	https://doi.org/10.1186/s12858-019-0115-3	article	more than 50 samples per group for both discovery and validation cohort	Case-control study	balanced error, accuracy, sensitivity, specificity, AUC (cross-validation, training + validation cohort)	cross-validation + external cohort validation

"We address gene selection and machine learning methods for cancer classification using microarray gene expression data. Due to the high dimensionality of microarray data, traditional gene selection algorithms are filter-based, focusing on intrinsic properties of the data such as distance, dependency, and correlation. These methods are fast but select far too many genes to use for the classification task. In this work, we present a novel hybrid filter-wrapper gene subset selection algorithm that is an improved modification of our prior algorithm. Our proposed method employs interaction information to rank candidate genes to add into a gene subset. [...] Experimental results on ten public cancer microarray data sets show that our method consistently outperforms prior gene selection algorithms in terms of classification accuracy, while requiring a small number of selected genes."

"Triple-negative breast cancer (TNBC) has attracted more attention compared with other breast cancer subtypes due to its aggressive nature, poor prognosis, and chemotherapy remains the mainstay of treatment with no other approved targeted therapy. Therefore, the study aimed to discover more promising therapeutic targets and investigating new insights of biological mechanism of TNBC. Six microarray data sets consisting of 463 non-TNBC and 405 TNBC samples were mined from Gene Expression Omnibus. [...] A naïve Bayes based classifier built using the expression profiles of 16 features (hub genes) accurately and reliably classify TNBC from non-TNBC samples in the validation test data set with a receiver operating curve of 0.93 to 0.98."

"We built a personalized prediction model based on clinical and genomic data that outperformed IPS5 and IPS5-R in predicting OS and AML transformation. The new model gives survival probabilities at different time points that are unique for a given pt. Incorporating clinical and mutational data outperformed a mutations only model even when cytogenetics and age were added"

"While treatment with the hypomethylating agents (HMA) azacitidine (AZA) and decitabine (DAC) improves cytogenetics and prolongs survival in MDS patients (pts), response is not guaranteed. Identification of non-responders could prevent prolonged exposure to ineffective therapy, avoid toxicities and decrease unnecessary costs. We developed an unbiased framework to study the association of several mutations in predicting response to HMA, analogous to Netflix or Amazon's recommender system in which customers who bought products A and B is likely to buy C; pts who have a mutation in gene A, and B are likely to respond or not respond to HMA. [...] When applying the model to the validation cohort, the sensitivity of these genomic biomarkers for no response to HMA was 1, specificity for response to HMA was 1, and accuracy was .85."

"Multiple proteomics could improve understanding and risk prediction of major adverse cardiovascular events (MACE) in type 2 diabetes. This study assessed 80 cardiovascular and inflammatory proteins for biomarker discovery and prediction of MACE in type 2 diabetes. [...] Addition of the B0-protein assay to the established risk model improved discrimination in the separate validation sample from 68.8 (95% CI, 68.2%-68.9%) to 74.8% (95% CI, 74.6%-75.1%)."

"Gene expression microarray data for 109 tumor cell lines with known sensitivity to the death ligand cytokine tumor necrosis factor-related apoptosis-inducing ligand (TRAIL) was used to identify genes with potential functional relationships determining responsiveness to TRAIL-induced apoptosis. The machine learning technique Random Forest in the statistical environment "R" with backward elimination was used to identify the key predictors of TRAIL. [...] Prediction accuracy was assessed by calculating the area under the receiver operator curve using an independent dataset. We show that the gene panel identified could predict TRAIL-sensitivity with a very high degree of sensitivity and specificity (AUC = 0.84). The genes in the panel are co-regulated and at least 40% of them functionally interact in signal transduction pathways that regulate cell death and cell survival, cellular differentiation and morphogenesis. Importantly, only 12% of the TRAIL predictor genes were differentially expressed highlighting the importance of functional interactions in predicting the biological responses."

"Prostate-specific antigen (PSA) is the most widely used serum biomarker for early detection of prostate cancer (PCA). Nevertheless, PSA level can be falsely elevated due to prostatic enlargement, inflammation or infection, which limits the PSA test specificity. The objective of this study is to use a machine learning approach for the analysis of mass spectrometry data to discover more reliable biomarkers that distinguish PCA from benign specimens. [...] From the new marker selection algorithm, a panel of 26 peaks achieved an accuracy of 80.7%, a sensitivity of 83.5%, a specificity of 78.4%, a positive predictive value (PPV) of 87.9%, and a negative predictive value (NPV) of 68.2%. On the other hand, when PSA alone was used (with a cutoff of 4.0 ng/ml), a sensitivity of 66.7%, a specificity of 53.6%, a PPV of 73.5%, and a NPV of 45.4% were obtained."

"A central challenge in systems biology and medical genetics is to understand how interactions among genetic loci contribute to complex phenotypic traits and human diseases. While most studies have so far relied on statistical modeling and association testing procedures, machine learning and predictive modeling approaches are increasingly being applied to mining genotype-phenotype relationships, also among those associations that do not necessarily meet statistical significance at the level of individual variants, yet still contributing to the combined predictive power at the level of variant panels. Network-based analysis of genetic variants and their interaction partners is another emerging trend by which to explore how sub-network level features contribute to complex disease processes and related phenotypes. In this review, we describe the basic concepts and algorithms behind machine learning-based genetic feature selection approaches, their potential benefits and limitations in genome-wide setting, and how physical or genetic interaction networks could be used as a priori information for providing improved predictive power and mechanistic insights into the disease networks."

"Chronic lymphocytic leukemia (CLL) is the most frequent lymphoproliferative syndrome in western countries. CLL evolution is frequently indolent, and treatment is mostly reserved for those patients with signs or symptoms of disease progression. In this work, we used RNA sequencing data from the International Cancer Genome Consortium CLL cohort to determine new gene expression patterns that correlate with clinical evolution. We determined that a 250-gene expression signature, in addition to immunoglobulin heavy chain variable region (IGHV) mutation status, stratifies patients into four groups with notably different time to first treatment. This finding was confirmed in an independent cohort. Similarly, we present a machine learning algorithm that predicts the need for treatment within the first 5 years following diagnosis using expression data from 2,198 genes. This predictor achieved 90% precision and 89% accuracy in predicting CLL cases."

"Late-Onset Alzheimer's Disease (LOAD) is a leading form of dementia. There is no effective cure for LOAD, leaving the treatment efforts to depend on preventive cognitive therapies, which stand to benefit from the timely estimation of the risk of developing the disease. We conducted systematic comparisons of representative Machine Learning models for predicting LOAD from genetic variation data provided by the Alzheimer's Disease Neuroimaging Initiative (ADNI) cohort. Our experimental results demonstrate that the classification performance of the best models tested yielded ~72% of area under the ROC curve"

Preprint BMJ Open 2021; 5(1):e003674. doi: <https://doi.org/10.1136/bmjopen-2020-025444>. Published online first April 18, 2024 by guest. Protected by copyright.

1										
2										
3										
4										
5										
6	Orienko, A and Moore, J H and Orzechowski, J and Ollon, R S and Cairns, J and Caraballo, P J and Weinshtaub, R M and Wang, L W and Breitenstein, M K	Considerations for automated machine learning in clinical metabolic profiling: Altered homocysteine on plasma concentration associated with metformin exposure	Pacific Symposium on Metabolism	2018	400-471	2018	USA	article	accuracy (training + test set split)	training + test set
7										
8	Pandey, G and Pandey, O P and Rogers, A J and Ahen, M E and Hoffman, G E and Raby, B A and Weiss, S T and Schadt, E E and Bunyavith, S	A Nasal Brush-based Classifier of Asthma Identified by Machine Learning Analysis of Nasal RNA Sequencing Data	Scientific Reports	8	15-15	2018	USA	article	AUC (5-fold CV)	cross-validation
9										
10										
11										
12										
13										
14										
15										
16										
17	154 Parker, B J and Gunter, S and Bedo, J	Stratification bias in low signal microarray studies	BMC Bioinformatics	8	326-326	2007	Australia	article	review (not applicable)	review
18										
19	Patil, S and Awan, K H and Arakati, G and Senviratne, C J and Muddur, N and Malik, S and Ferrar, M and Rahimi, S and Brennan, P A	Machine learning and its potential applications to the genomic study of head and neck cancer: A systematic review	Journal of Oral Pathology & Medicine	48	9	773-779	2019	India	article	review (not applicable)
20										
21										
22										
23										
24	Porter, D and Goodyear, C S and Nijjar, J S and Messow, M and Seber, S and Mudaliar, A and McInnes, I B	Predicting the response to TNF inhibition or B cell depletion therapy from peripheral whole blood gene expression profiles in patients with rheumatoid arthritis	Arthritis and Rheumatology	68	4130-4131	2016		meeting abstract	Cases only (drug response study)	sensitivity, specificity, PPV and NPV (10 fold CV)
25										
26										
27										
28										
29	157 Pratt, A G and Swan, D and Richardson, S and Wilson, G and Hilkens, C and Young, D and Isaacs, J D	"Microarray analysis of 111 RNA samples was performed. [...] Machine learning approaches were used to test the utility of a classification model amongst an independent validation cohort of 62 patients presenting with UA." "A high-throughput DNA methylation dataset (100 samples) of ESCC from the Cancer Genome Atlas (TCGA) project was analyzed and validated along with another independent dataset (12 samples) from the Gene Expression Omnibus (GEO) database. [...] The candidate CpG sites as well as their adjacent regions were further validated in 94 pairs of ESCC tumor and adjacent normal tissues from the Chinese Han population using the targeted bisulfite sequencing method."	Arthritis and Rheumatism	m	63	10	2011	article	Case-control study	sensitivity, specificity training and validation cohort)
30										
31										
32										
33										
34	Pu, W and Wang, C and Chen, S and Zhao, D and Zhou, Y and Ma, Y and Yang, W and Li, C and Huang, S and Jin, L and Guo, S and Wang, J and Wang, M	Targeted bisulfite sequencing identified a panel of DNA methylation-based biomarkers for esophageal squamous cell carcinoma (ESCC)	Clin Epigenetics	9	129-129	2017	China	article	AUC, accuracy, sensitivity, specificity (5-fold cross-validation + test set)	cross-validation + test set
35										
36										
37										
38										
39	Pujos-Guillot, E and Bertrand, J and Camilleau, M and Piérat, M and Brandolini, M and Fernandes, A and Matta, J and Lévy-Marchal, C and Renouard, S and Comte, B	Contribution of an integrative multi-omic approach in the metabolic syndrome prediction: A nested case-control study	Drug Metabolism and Personalized Therapy	31	4	e433-434	2016	meeting abstract	n=103 biom score n=70 biom adequate for gestational age	Case-control study error rate (training/test split + validation set)
40										
41										
42										
43										
44										
45										
46										

"With the maturation of metabolomics science and proliferation of biobanks, clinical metabolic profiling is an increasingly opportunistic frontier for advancing translational clinical research. Automated Machine Learning (AutoML) approaches provide exciting opportunity to guide feature selection in agnostic metabolic profiling endeavors, where potentially thousands of independent data points must be evaluated. In previous research, AutoML, using high-dimensional data of varying types has been demonstrated to outperform traditional approaches. However, several considerations for application in clinical metabolic profiling remain to be evaluated. Particularly, regarding the robustness of AutoML to identify and adjust for common clinical confounders. In this study, we present a focused case study regarding AutoML considerations for using the Tree-Based Optimization Tool (TBO) in metabolic profiling of exposure to metformin in a biobank cohort. [...] First, we propose a tandem rank accuracy measure to guide agnostic feature selection and corresponding threshold determination in clinical metabolic profiling endeavors. Second, while AutoML, using default parameters, demonstrated potential to lack sensitivity to detect confounding clinical covariates, we demonstrated residual training and adjustment of metabolic features as an easily applicable approach to ensure AutoML adjustment for potential confounding characteristics. Finally, we present increased homogeneity with long term exposure to metformin as a potentially novel, non-replicated metabolic association suggested by TBO as an association not identified in parallel clinical metabolic profiling endeavors." "Asthma is a common, under-diagnosed disease affecting all ages. We sought to identify a nasal brush-based classifier of mild/moderate asthma. 150 subjects with mild/moderate asthma and controls underwent nasal brushing and RNA sequencing of nasal samples. A machine learning-based pipeline identified an asthma classifier consisting of 30 genes interpreted via an L2-regularized logistic regression classification model. This classifier performed with strong predictive value and sensitivity across eight test sets." "When analyzing microarray and other small sample size biological datasets, care is needed to avoid various biases. We analyze a form of bias, stratification bias, that can substantially affect analyses using sample-reuse validation techniques and lead to inaccurate results. This bias is due to imperfect stratification of samples in the training and test sets and the dependency between these stratification errors, i.e. the variations in proportions in the training and test sets are negatively correlated. We show that when estimating the performance of classifiers on low signal datasets (i.e. those which are difficult to classify), which are typical of many prognostic microarray studies, commonly used performance measures can suffer from a substantial negative bias. For error rate this bias is only severe in quiet restricted situations, but can be much larger and more frequent when using ranking measures such as the receiver operating characteristic (ROC) curve and area under the ROC (AUC). [...] The classification error rate can have large negative biases for balanced datasets, whereas the AUC shows substantial pessimistic biases even for imbalanced datasets. [...] Stratification bias can substantially affect several performance measures. In computing the AUC, the strategy of pooling the test samples from the various folds of cross-validation can lead to large biases; computing it as the average of per-fold estimates avoids this bias and is thus the recommended approach. As a more general solution applicable to other performance measures, we show that stratified repeated holdout and a modified version of k-fold cross-validation, balanced, stratified cross-validation and balanced leave-one-out cross-validation, avoids this bias. Therefore for model selection and evaluation of microarray and other small biological datasets, these methods should be used and unstratified versions avoided. In particular, the commonly used (unbalanced) leave-one-out cross-validation should not be used to estimate AUC for small datasets." "The aim of this systematic review was to evaluate the existing literature and assess the application of machine learning of genomic data in head and neck cancer (HNC). [...] Two studies each evaluated oral cancer and laryngeal cancer, while other one study each evaluated esophageal cancer and oropharyngeal cancer. The majority of studies employed support vector machine (SVM) as a ML technique. Among the included studies, the accuracy rates for ML techniques ranged from 56.7% to 99.4%. Our findings showed that ML techniques for the analysis of genomic data can play a role in the prognostic prediction of HNC." "The ORBIT study demonstrated that rituximab is non-inferior to a TNFi-first strategy in biologic naïve, sero-positive patients with active rheumatoid arthritis (RA) over 12 months [Lancet doi.org/10.1016/S0140-6736(16)03809-9]. However, a significant proportion of patients failed to respond to their first biologic drug and switched to an alternative. The ability to identify and stratify these patients prior to treatment would improve patient care and optimize the use of scarce financial resources. The aim of this study was to identify peripheral blood transcriptional biomarkers in the ORBIT cohort that can predict subsequent response/non-response to biologic therapy. Three gene sets were identified using support vector machine (SVM) recursive feature elimination. These predicted: general responsiveness to both TNFi and rituximab therapy (8 genes), response to TNFi therapy (23 genes) or rituximab (23 genes) respectively. When tested on the validation set, these models resulted in ROC plots with an AUC of 91.6% for general responsiveness, 89.7% for TNFi responsiveness, and 85.7% for rituximab response." "The diagnosis of seronegative rheumatoid arthritis (RA) remains challenging in the early arthritis clinic. Recent GWAS data strongly implicate CD4+ T cells in the pathogenesis of seropositive RA. Our objectives were to identify biomarker(s) present in CD4+ T cells, or in serum, that identified patients with undifferentiated arthritis (UA) destined to develop seronegative RA. [...] Machine learning approaches were used to test the utility of a classification model amongst an independent validation cohort of 62 patients presenting with UA. [...] A 12-gene expression "signature" predicted the subsequent development of RA amongst ACRA-negative UA patients in the validation cohort (sensitivity 85%, specificity 75%). The signature had a predictive value equivalent to the Leiden score in these patients and provided enhanced predictive power in combination with the Leiden score. The 12-gene signature confirmed an over-representation of STAT3 target genes, and pathway analysis confirmed that genes functionally involved with CD4+ T cell survival, including STAT pathway components, were downregulated in early RA." "DNA methylation has been implicated as a promising biomarker for precise cancer diagnosis. However, limited DNA methylation-based biomarkers have been described in esophageal squamous cell carcinoma (ESCC). [...] A high-throughput DNA methylation dataset (100 samples) of ESCC from The Cancer Genome Atlas (TCGA) project was analyzed and validated along with another independent dataset (12 samples) from the Gene Expression Omnibus (GEO) database. The methylation status of peripheral blood mononuclear cells and peripheral blood leukocytes from healthy controls was also utilized for biomarker selection. The candidate CpG sites as well as their adjacent regions were further validated in 94 pairs of ESCC tumor and adjacent normal tissues from the Chinese Han population using the targeted bisulfite sequencing method. [...] Eight statistical models along with five-fold cross-validation were further applied, in which the SVM model reached the best accuracy in both training and test dataset (accuracy = 0.82 and 0.80, respectively)." "The rising worldwide prevalence of metabolic syndrome (MetS), a cluster of cardiometabolic risk factors of predictive of type 2 diabetes, relates largely to increasing obesity and sedentary but also to early metabolic life events [1]. Objective: The objective of the study was to identify predictive biomarkers of evolution toward MetS 8 years later, and to bring new knowledge about this pathological state using a multidisciplinary approach in an at-risk population (subjects with small birth weight). [...] Individual predictive models were first built using linear logistic regressions from the omics datasets. Metabolic and proteomic data were finally integrated using random forests to determine whether multidimensional models improve prediction. The resulting models based on either 4 metabolites or 4 proteins showed good performance: 22% misclassification on training set, 25% on validation set vs 11% misclassification on training set, 33% on validation set, respectively. Multi-omic data integration improved performance: 15% misclassification on training set, 12% misclassification on training set, 8% on validation set."

"With the maturation of metabolomics science and proliferation of biobanks, clinical metabolic profiling is an increasingly opportunistic frontier for advancing translational clinical research. Automated Machine Learning (AutoML) approaches provide exciting opportunity to guide feature selection in agnostic metabolic profiling endeavors, where potentially thousands of independent data points must be evaluated. In previous research, AutoML, using high-dimensional data of varying types has been demonstrated to outperform traditional approaches. However, several considerations for application in clinical metabolic profiling remain to be evaluated. Particularly, regarding the robustness of AutoML to identify and adjust for common clinical confounders. In this study, we present a focused case study regarding AutoML considerations for using the Tree-Based Optimization Tool (TBO) in metabolic profiling of exposure to metformin in a biobank cohort. [...] First, we propose a tandem rank accuracy measure to guide agnostic feature selection and corresponding threshold determination in clinical metabolic profiling endeavors. Second, while AutoML, using default parameters, demonstrated potential to lack sensitivity to low-effect confounding clinical covariates, we demonstrated residual training and adjustment of metabolic features as an easily applicable approach to ensure AutoML adjustment for potential confounding characteristics. Finally, we present increased homogeneity with long term exposure to metformin as a potentially novel, non-replicated metabolic association suggested by TBO as an association not identified in parallel clinical metabolic profiling endeavors." "Asthma is a common, under-diagnosed disease affecting all ages. We sought to identify a nasal brush-based classifier of mild/moderate asthma. 150 subjects with mild/moderate asthma and controls underwent nasal brushing and RNA sequencing of nasal samples. A machine learning-based pipeline identified an asthma classifier consisting of 30 genes interpreted via an L2-regularized logistic regression classification model. This classifier performed with strong predictive value and sensitivity across eight test sets." "When analyzing microarray and other small sample size biological datasets, care is needed to avoid various biases. We analyze a form of bias, stratification bias, that can substantially affect analyses using sample-reuse validation techniques and lead to inaccurate results. This bias is due to imperfect stratification of samples in the training and test sets and the dependency between these stratification errors, i.e. the variations in proportions in the training and test sets are negatively correlated. We show that when estimating the performance of classifiers on low signal datasets (i.e. those which are difficult to classify), which are typical of many prognostic microarray studies, commonly used performance measures can suffer from a substantial negative bias. For error rate this bias is only severe in quiet restricted situations, but can be much larger and more frequent when using ranking measures such as the receiver operating characteristic (ROC) curve and area under the ROC (AUC). [...] The classification error rate can have large negative biases for balanced datasets, whereas the AUC shows substantial pessimistic biases even for imbalanced datasets. [...] Stratification bias can substantially affect several performance measures. In computing the AUC, the strategy of pooling the test samples from the various folds of cross-validation can lead to large biases; computing it as the average of per-fold estimates avoids this bias and is thus the recommended approach. As a more general solution applicable to other performance measures, we show that stratified repeated holdout and a modified version of k-fold cross-validation, balanced, stratified cross-validation and balanced leave-one-out cross-validation, avoids this bias. Therefore for model selection and evaluation of microarray and other small biological datasets, these methods should be used and unstratified versions avoided. In particular, the commonly used (unbalanced) leave-one-out cross-validation should not be used to estimate AUC for small datasets." "The aim of this systematic review was to evaluate the existing literature and assess the application of machine learning of genomic data in head and neck cancer (HNC). [...] Two studies each evaluated oral cancer and laryngeal cancer, while other one study each evaluated esophageal cancer and oropharyngeal cancer. The majority of studies employed support vector machine (SVM) as a ML technique. Among the included studies, the accuracy rates for ML techniques ranged from 56.7% to 99.4%. Our findings showed that ML techniques for the analysis of genomic data can play a role in the prognostic prediction of HNC." "The ORBIT study demonstrated that rituximab is non-inferior to a TNFi-first strategy in biologic naïve, sero-positive patients with active rheumatoid arthritis (RA) over 12 months [Lancet doi.org/10.1016/S0140-6736(16)03809-9]. However, a significant proportion of patients failed to respond to their first biologic drug and switched to an alternative. The ability to identify and stratify these patients prior to treatment would improve patient care and optimize the use of scarce financial resources. The aim of this study was to identify peripheral blood transcriptional biomarkers in the ORBIT cohort that can predict subsequent response/non-response to biologic therapy. Three gene sets were identified using support vector machine (SVM) recursive feature elimination. These predicted: general responsiveness to both TNFi and rituximab therapy (8 genes), response to TNFi therapy (23 genes) or rituximab (23 genes) respectively. When tested on the validation set, these models resulted in ROC plots with an AUC of 91.6% for general responsiveness, 89.7% for TNFi responsiveness, and 85.7% for rituximab response." "The diagnosis of seronegative rheumatoid arthritis (RA) remains challenging in the early arthritis clinic. Recent GWAS data strongly implicate CD4+ T cells in the pathogenesis of seropositive RA. Our objectives were to identify biomarker(s) present in CD4+ T cells, or in serum, that identified patients with undifferentiated arthritis (UA) destined to develop seronegative RA. [...] Machine learning approaches were used to test the utility of a classification model amongst an independent validation cohort of 62 patients presenting with UA. [...] A 12-gene expression "signature" predicted the subsequent development of RA amongst ACRA-negative UA patients in the validation cohort (sensitivity 85%, specificity 75%). The signature had a predictive value equivalent to the Leiden score in these patients and provided enhanced predictive power in combination with the Leiden score. The 12-gene signature confirmed an over-representation of STAT3 target genes, and pathway analysis confirmed that genes functionally involved with CD4+ T cell survival, including STAT pathway components, were downregulated in early RA." "DNA methylation has been implicated as a promising biomarker for precise cancer diagnosis. However, limited DNA methylation-based biomarkers have been described in esophageal squamous cell carcinoma (ESCC). [...] A high-throughput DNA methylation dataset (100 samples) of ESCC from The Cancer Genome Atlas (TCGA) project was analyzed and validated along with another independent dataset (12 samples) from the Gene Expression Omnibus (GEO) database. The methylation status of peripheral blood mononuclear cells and peripheral blood leukocytes from healthy controls was also utilized for biomarker selection. The candidate CpG sites as well as their adjacent regions were further validated in 94 pairs of ESCC tumor and adjacent normal tissues from the Chinese Han population using the targeted bisulfite sequencing method. [...] Eight statistical models along with five-fold cross-validation were further applied, in which the SVM model reached the best accuracy in both training and test dataset (accuracy = 0.82 and 0.80, respectively)." "The rising worldwide prevalence of metabolic syndrome (MetS), a cluster of cardiometabolic risk factors of predictive of type 2 diabetes, relates largely to increasing obesity and sedentary but also to early metabolic life events [1]. Objective: The objective of the study was to identify predictive biomarkers of evolution toward MetS 8 years later, and to bring new knowledge about this pathological state using a multidisciplinary approach in an at-risk population (subjects with small birth weight). [...] Individual predictive models were first built using linear logistic regressions from the omics datasets. Metabolic and proteomic data were finally integrated using random forests to determine whether multidimensional models improve prediction. The resulting models based on either 4 metabolites or 4 proteins showed good performance: 22% misclassification on training set, 25% on validation set vs 11% misclassification on training set, 33% on validation set, respectively. Multi-omic data integration improved performance: 15% misclassification on training set, 12% misclassification on training set, 8% on validation set."

BMJ Open 2021; 15:e003674. doi:10.1136/bmjopen-2021-003674

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46

160	Puzstall, L and Hess, K R	Clinical trial design for microarray predictive marker discovery and assessment	Ann Oncol	15	12	1791-1797	2004	USA	http://dx.doi.org/10.1093/annonc/mdh464	article	review (not applicable)	review				
161	Rao, R and Dean, K and Migasawa, B and Somwarshi, P and Doyle, F	Validation of a Multi Omic Biomarker Panel and Analysis of Disease Progression Trajectories in a Novel Longitudinal PTSD Cohort	Biological Psychiatry	85	10	596-596	2019		http://dx.doi.org/10.1016/j.biopsych.2019.03.042	meeting abstract	83 PTSD positive cases and 83 PTSD negative matched controls, and subsequently refined and validated in a cohort of 29 PTSD cases and 40 controls	Case-control study	AUC, accuracy, sensitivity, specificity (training / test set + external validation)	validation	external cohort validation	"This manuscript reviews methodological and statistical issues relevant to clinical trial design to discover and validate multiple predictors of response to therapy." "Signals associated with PTSD development might emerge across multiple levels of physiological function. Diagnostic classifiers synthesizing signals from several single-layer molecular features into a multi-omic panel can improve diagnostic performance compared to any individual molecular signature. [...] Single and multi-omic classifiers were initially identified in a cohort of 83 PTSD positive cases and 83 PTSD negative matched controls, and subsequently refined and validated in a cohort of 29 PTSD cases and 40 controls. A novel longitudinal cohort of 1800 active duty soldiers is used for external validation [...] We previously found that the multi-omic panel results in a small improvement in diagnostic performance in comparison to individual single-omic panels in the initial training and validation cohorts (AUC=0.80, 77% accuracy, sensitivity, 73% specificity). Preliminary external validation in the longitudinal cohort suggests that single-omic, metabolic panels constituting the multi-omic panel are significantly associated with PTSD status."
162	Rappoport, M and Shamir, R	Multi-omic and multi-view clustering algorithms: Review and cancer benchmark	Nucleic Acids Research	46	20	10546-10562	2018	Israel	http://dx.doi.org/10.1093/nar/nky121	article	review (not applicable)	review				
163	Revee, J and Madhl-Thomsen, K S and Halloran, P F	Using ensembles of machine learning classifiers to maximize the accuracy and stability of molecular biopsy interpretation	American Journal of Transplantation	19		452-463	2019		http://dx.doi.org/10.1111/ajt.15490	meeting abstract	"1679 kidney transplant biopsies were repeatedly split at random into two training sets (N=600 each) and a test set (N=179)"	Case-control study	accuracy (training and validation set)	training + test set		
164	Revee, J and Sellares, J and De Freitas, D and Ginecke, G and Bromberg, J and Matas, A and Halloran, P	Predicting graft failure in kidney transplant patients: A combined clinical/molecular approach to analyzing biopsies	American Journal of Transplantation	13		109-109	2013		http://dx.doi.org/10.1111/ajt.12295	meeting abstract	562 patients (1 biopsy per patient)	Case-control study	accuracy (training / test set split)	training + test set		
165	Reinhold, W C and Varma, S and Rajapakse, V N and Luna, A and Sousa, F and Koh, K W and Pommer, Y G	Using drug response data to identify molecular effectors and metabolic "omic" data to identify candidate drugs in cancer	Hum Genet	134	1	3-11	2015	USA	http://dx.doi.org/10.1007/s00438-014-0737-6	article	review (not applicable)	review				
166	Resom, H W and Varghese, R S and Abdel-Hamid, M and Eissa, S A and Saha, D and Goldman, L and Petrosino, E F and Canden, T P and Venstra, T D and Loffredo, C A and Goldman, R	Analysis of mass spectral serum profiles for biomarker selection	Bioinformatics	21	21	4039-4045	2005	USA	http://dx.doi.org/10.1093/bioinformatics/bti070	article	411 sera samples (199 from HCC patients and 212 from matched healthy individuals)	Case-control study	accuracy, sensitivity, specificity (fold cross-validation and bootstrapping methods)	cross-validation		
167	Ritari, J and Hyvrikanen, K and Kukola, S and Tuho-Renno, M and Nittynuopio, P and Nihminen, A and Salonen, U and Puskonen, M and Vaini, L and Kwan, T and Pastinen, T and Partanen, J	Genomic prediction of relapse in recipients of allogeneic haematopoietic stem cell transplantation	Leukemia	33	1	240-248	2019	Canada	http://dx.doi.org/10.1038/s41375-018-0223-3	article	"We studied 151 graft recipients with HLA-matched sibling donors by sequencing the whole-exome, active immunoregulatory regions, and the full MHC region."	Case only (relapse prediction)	AUC with confidence intervals (LOOCV)	cross-validation		
168	Roder, J and Oliveira, C and Net, L and Tsipin, M and Lindtli, B and Roder, H	A dropout-regulated classifier development approach optimized for oncotic data	Bmc Bioinformatics	20		14-14	2019	USA	http://dx.doi.org/10.1186/s12859-019-2022-2	article	"For the GSE50081 cohort expression profiling was performed on RNA from frozen, resected tumor tissue from 141 subjects with stage I or II NSCLC [...] Expression profiling for the GSE42127 cohort was performed for 176 subjects with stage I-IV NSCLC."	Case-control study	AUC (training and validation cohort)	external cohort validation		

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46

Author(s)	Title	Year	Country	Study Type	Key Findings	Validation
169 Rodriguez Ortiz, M E and Perrella, C and Rodriguez, M and Zurbig, P and Michals, H and Ortiz, A Rosenberg, S and Durcay, F and Alentorn, C and Delhais, C and Elaroui, N and Kamoun, A and Mané, Y and Tangou, M L and De Reynies, A and Mokhtari, K and Figarella-Branger, D and Delattre, J Y and Idbahl, A and Adam, C and Andrad, M and Aubriot-Lorton, M H and Bauche, L and Beauchêne, P and Blehert, L and Blechet, C and Campone, M and Carpentier, A and Carvajal, L and Catal-Hatem, D and Lhermitte, B and Chiforeanu, D and Chinot, O and Cohen-Moyal, E and Colin, P and Cruet, T and Dam-Hieu, P and Desseaux, C and Desse, N and Dhermain, F and Diebold, M D and Limier, S and Fallois, T and Fesneau, M and Fontaine, D and Gallard, S and Forest, F and Gauthier, G and Gaudier, C and Ghiringhelli, F and Godfrand, C and Guye, E M and Elouadhani-Hamdi, S and Honorat, J and Khalil, T and Labrousse, F and Lahiani, W and Langlois, O and Laquerrière, A and Larrieu-Ciron, D and Lechapt-Zakman, E and Lousseau, H and Lopez, S and Loussouart, D and Mouraux, C A and Mene, P and Mihai, M I and Milin, S and Fotso, M J M and Noel, G and Parker, F and Pett, A and Quintin-Roue, I and Ramirez, C and Rousseau, A and Rousselet-Denis, C and Ricard, D and Richard, P and Rigau, V and Roussel, G and Svestov, H and Tortat, M C and Vandenberg, O and Vauleon, E and Villa, C and Zemmoura, I and Desseaux, C and Svestov, H and Mene, P and Rousseau, A and Cruet, T and Lopez, S and Mihai, M I and Pett, A and Adam, C and Parker, F and Carpentier, A and Dam-Hieu, P and Quintin-Roue, I and Eimer, S and Lousseau, H and Blehert, L and Lechapt-Zakman, E and Godfrand, C and Khalil, T and Catal-Hatem, D and Fallois, T and Andrad, M and Carvajal, L and Richard, P and Lahiani, W and Aubriot-Lorton, H and Ghiringhelli, F and Mouraux, C A and Ramirez, C and Guye, E M and Labrousse, F and Durcay, F and Jouve, A and Figarella-Branger, D and Chinot, O and Bauche, L and Rigau, V and Beauchêne, P and Gallard, G and Campone, M and Loussouart, D and Fontaine, D and Vandenberg, O and Desse, N and Fesneau, M and Delhais, C and Delattre, J Y and Elouadhani-Hamdi, S and Ricard, D and Larrieu-Ciron, D and Milin, S and Colin, P and Diebold, M D and Chiforeanu, D and Vauleon, E and Langlois, O and Laquerrière, A and Forest, F and Mene, P and Andrad, M and Roussel, G and Lhermitte, B and Noel, G and Gallard, S and Villa, C and Desse, N and Cohen-Moyal, E and Uro-Coste, E and Dhermain, F and Network, Pola	Novel Urinary Biomarkers For Improved Prediction Of Prognostic Risk In Early Chronic Disease Stages And In High Risk Individuals Without Chronic Kidney Disease	2018	Germany	Rapid progressors (n = 342) Non-rapid progressors (n = 1140)	AUC, sensitivity, specificity, NPV, PPV (LOOCV + external validation cohort)	cross-validation + external cohort validation
171 Rychkov, D and Sirota, M and Lin, C	Leveraging publicly available gene expression data and applying machine learning to identify novel biomarkers for rheumatoid arthritis	2018		Meeting abstract	Cohen's kappa, sensitivity, specificity (5-fold CV)	cross-validation
172 Sahli, G and Mittal, K and Rida, P and Janssen, E A M and Gognigni, K and Ajeza, R	Panoptic view of prognostic models for personalized breast cancer management	2019	USA	Review (not applicable)	Review	Review (not applicable)
173 Schreiber, T and Lustrek, M and Schmidt, R and Repolster, D and Fuxlen, G	Tissue-based Alzheimer gene expression markers comparison of Alzheimer's disease with other neurodegenerative BMC and investigation of redundancy in small biomarker sets	2012	Germany	Article	PLURI and AD dataset contain > 50 samples per group	cross-validation
174 Shaif, A and Nguyen, T and Peyvandpour, A and Nguyen, H and Draglicki, S	A multi-cohort and multi-omics meta-analysis framework to identify network-based gene signatures	2019	United States	Article	622 samples; 533 samples from GBM patients and 89 from healthy (non-tumor) individuals	external cohort validation
175 Shaik, A and Acharya, C and Smetzer, S and Iyerli, H K and Acharya, K S	Non-invasive diagnosis of endometriosis using machine learning instead of the operating room	2019		Article	derived from spectral decomposition of the discovery dataset (n: 148) to predict the presence of endometriosis*	cross-validation + external cohort validation
176 Shan, L and Chen, Y A and Davis, L and Han, G and Zhu, W and Molina, D A and Arango, H and LaPolla, J P and Hoffman, M S and Sellers, T and Kirby, D and Nicotia, S V and Suthphen, R	Measurement of phospholipids may improve diagnostic accuracy in ovarian cancer	2012	USA	Article	*total of 1057 women with suspected ovarian cancer were enrolled [...] Only patients who underwent surgery based on clinical suspicion of ovarian cancer were eligible and if a patient was diagnosed with EOC, surgical staging was documented (including 233 in whom EOC was confirmed [...]) A total of 211 cases and 212 benigns was included in the analysis.*	cross-validation
177 and Tsui, K H and Juo, C G and Wu, K P	Metabolic profile discovery for the detection of bladder cancer by comparative metabolomics	2017	China	Article	*metabolic profiles of 87 samples from bladder cancer patients and 65 samples from hermia patients*	cross-validation + test set

"Chronic kidney disease is associated with increased risk of CKD progression and death. Therapeutic approaches to limit progression are limited. Developing tools for the early identification of those individuals most likely to progress will allow enrolling clinical trials in high risk early CKD patients. The CKD273 classifier is a panel of 273 urinary peptides that enables early detection of CKD and prognosis of progression. We have generated urine capillary electrophoresis-mass spectrometry-based peptideomics CKD273 subclassifiers specific for CKD stages to allow the early identification of patients at high risk of CKD progression. [...] In individuals with eGFR > 60 mL/min/1.73 m² and albuminuria <30 mg/day, the CKD273 subclassifiers predicted rapid eGFR loss with AUC ranging from 0.79 (0.743-0.844) to 0.736 (0.689-0.780). The association between CKD273 subclassifiers and rapid progression remained significant after adjustment for age, sex, albuminuria, DM, baseline eGFR, and systolic blood pressure."

"1p19q codeleted anaplastic gliomas have variable clinical behavior. We have recently shown that the common Sp12.3 allelic loss is an independent prognostic factor in this tumor type. The aim of this study is to identify less frequent genomic copy number variations (CNVs) with clinical importance that may shed light on molecular oncogenesis of this tumor type. [...] Computational biology and feature selection based on the random forests method were used to identify CNV events associated with overall survival and other clinical pathological variables. [...] Several recurrent CNV events, detected in anaplastic oligodendrogliomas, enable better survival prediction. More importantly, they help in identifying potential gene targets for understanding oncogenesis and for personalized therapy."

"Diagnosis and monitoring the disease progression of RA is challenging requiring a combination of imaging techniques and blood tests. There is currently no biochemical test for detection of early-stage disease. In this study, we aimed to define a Rheumatoid Arthritis meta-profile and identify biomarkers by leveraging publicly available gene expression data with machine learning approaches. [...] Finally, we built a Random Forest classification model on the synovium data with these 5 genes. We applied 5-fold cross-validation with 10 repeats technique and used Cohen's Kappa statistic as a metric. We obtained Kappa equal 0.51 with sensitivity 0.86 and specificity 0.53 on the testing set. In the final step, we validated the prediction model on the whole blood data, resulting kappa of 0.57 with sensitivity 0.54 and specificity 0.86."

"The efforts to personalize treatment for patients with breast cancer have led to a focus on the characterization of genotypic and phenotypic heterogeneity among breast cancers. [...] This review summarizes the prognostic and predictive insights provided by commercially available gene expression based tests and other multivariate clinical-omics based prognostic/predictive models currently under development, and proposes a more inclusive multiparameter approach to tackling the challenging heterogeneity of breast cancer to individualize management."

"Alzheimer's disease has been known for more than 100 years and the underlying molecular mechanisms are not yet completely understood. The identification of genes involved in the processes in Alzheimer affected brain is an important step towards such an understanding. [...] Based on microarray data we identify potential biomarkers as well as biomarker combinations using three feature selection methods: information gain, mean decrease accuracy of random forest and a wrapper of genetic algorithm and support vector machine (GA/SVM). [...] Compared to the other methods, GA/SVM has the advantage of finding small, less redundant sets of genes that, in combination, show superior classification characteristics."

"Although massive amounts of condition-specific molecular profiles are being accumulated in public repositories every day, meaningful interpretation of these data remains a major challenge. In an effort to identify the biomarkers that describe the key biological phenomena for a given condition, several approaches have been developed over the past few years. However, the majority of these approaches either (i) do not consider the known intermolecular interactions, or (ii) do not integrate molecular data of multiple types (e.g. genomic, transcriptomic, proteomic, epigenomic), and thus potentially fail to capture the true biological changes responsible for complex diseases (e.g. cancer). In addition, these approaches often ignore the heterogeneity and study bias present in independent molecular cohorts. In this manuscript, we propose a novel multi-cohort and multi-omics meta-analysis framework that overcomes all three limitations mentioned above in order to identify robust molecular subnetworks that capture the key dynamic nature of a given biological condition. [...] We demonstrate the proposed framework by constructing subnetworks related to two complex diseases: glioblastoma and low-grade gliomas. We validated the identified subnetworks by showing their ability to predict patients' clinical outcome on multiple independent validation cohorts."

"Endometriosis affects an estimated 1 in 10 women during their reproductive years, and up to 30% to 50% of women with endometriosis may experience infertility. [...] A previous study developed classifiers for prediction of endometriosis in a cycle-phase specific manner by using margin tree classification within one dataset. Our aim was to build on this research by utilizing machine learning to predict and independently validate the presence or absence of endometriosis, regardless of cycle phase and other uterine pathology, through endometrial biopsy (EMB) samples. [...] We identified a 280 gene predictor of endometriosis using Random Forests that was found to predict the presence of endometriosis, regardless of the endometrial phase and other pathology, with an accuracy of 84% (area under ROC: 1.0 0.86; p-value: 6.14e-05), with a negative predictive value of 86% and a positive predictive value of 81%. We reduced model over-fitting by performing 10-fold cross-validation of our discovery data."

"More than two-thirds of women who undergo surgery for suspected ovarian neoplasm do not have cancer. Our previous results suggest phospholipids as potential biomarkers of ovarian cancer. In this study, we measured the serum levels of multiple phospholipids among women undergoing surgery for suspected ovarian cancer to identify biomarkers that better predict whether an ovarian mass is malignant. [...] The HE-SVM model using the measurements of specific combinations of phospholipids supplements clinical CA125 measurement and improves diagnostic accuracy. Specifically, the measurement of phospholipids improved sensitivity (identification of cases with preoperative CA125 levels below 35) among two types of cases in which CA125 performance is historically poor - early stage cases and those of mucinous histology. Measurement of phospholipids improved the identification of early stage cases from 68% (based on CA125) to 82%, and mucinous cases from 44% to 88%."

"In this study, we applied ultra-performance liquid chromatography time-of-flight mass spectrometry to profile metabolite profiles of 87 samples from bladder cancer patients and 65 samples from hermia patients. An OPLS-DA classification revealed that bladder cancer samples can be discriminated from hermia samples based on the profiles. A marker discovery pipeline selected six putative markers from the metabolomic profiles. [...] A machine learning model, decision trees, was built based on the metabolomic profiles and the six marker candidates. The decision tree obtained an accuracy of 76.60%, a sensitivity of 71.88%, and a specificity of 86.67% from an independent test."

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46

178	Sharma, A and Rani, R	C-HMOSHSA: Gene selection for cancer classification using multi-objective meta-heuristic and machine learning methods	Comput Methods Biomed	178	219-235	2019	India	http://dx.doi.org/10.1016/j.comp.2019.08.002	the proposed machine learning approaches tested on 7 microarray datasets, including a datasets with > 50 samples per group	Case-control study	accuracy (LOOCV + test set)	cross-validation + test set	"We have proposed a framework (C-HMOSHSA) for gene selection using multi-objective spotted hyena optimizer (MOSH) and salp swarm algorithm (SSA). The real life optimization problems with more than one objective usually face the challenge to maintain convergence and diversity. Salp Swarm Algorithm (SSA) maintains diversity but, suffers from the overhead of main-training the necessary information. On the other hand, the calculation of MOSHO requires low computational efforts hence is used for maintaining the necessary information. Therefore, the proposed algorithm is a hybrid algorithm that utilizes the features of both SSA and MOSHO to facilitate its exploration and exploitation capability. [...] Four different classifiers are trained on seven high-dimensional datasets using a subset of features (genes), which are obtained after applying the proposed hybrid gene selection algorithm. The results show that the proposed technique significantly outperforms existing state-of-the-art techniques." "The use of penalized logistic regression for cancer classification using microarray expression data is presented. Two dimension reduction methods are respectively combined with the penalized logistic regression so that both the classification accuracy and computational speed are enhanced. Two other machine learning methods, support vector machines and least-squares regression, have been chosen for comparison. It is shown that our methods have achieved at least equal or better results. They also have the advantage that the output probability can be explicitly given and the regression coefficients are easier to interpret."
179	Shen, Land Tan, C C	Dimension reduction-based penalized logistic regression for cancer classification using microarray data	IEEE/ACM Trans Comput Biol Bioinform	2	2	166-175	2005	Singapore	"the proposed machine learning approaches tested on 7 microarray datasets, including a datasets with > 50 samples per group	Case-control study	mean error + standard deviation (LOOCV, test set)	cross-validation + test set	"The use of penalized logistic regression for cancer classification using microarray expression data is presented. Two dimension reduction methods are respectively combined with the penalized logistic regression so that both the classification accuracy and computational speed are enhanced. Two other machine learning methods, support vector machines and least-squares regression, have been chosen for comparison. It is shown that our methods have achieved at least equal or better results. They also have the advantage that the output probability can be explicitly given and the regression coefficients are easier to interpret."
	Sherman, S I and Pagan, M and Huang, J and Lin, B and Diggins, J and Tom, E and Haugen, B and Tuttle, R M and Kennedy, G	Augmenting pre-operative risk of recurrence stratification in differentiated thyroid carcinoma using machine learning and high dimensional transcriptional data from thyroid FNA	Journal of Clinical Oncology	33	15		2015		"81 samples preoperatively collected in a previous study and post-surgically diagnosed as PTC [...] Each patient was categorized as either ATA low risk or ATA intermediate/high risk using established guidelines for recurrence risk stratification." (< 50 samples per group)	Cases only (risk of recurrence prediction)	AUC (cross-validation)	cross-validation	"Transcriptional data from FNA of thyroid nodules may improve the pre-operative prediction of post-operative recurrence. If independently validated in a sufficiently large number of patients, such molecular classifiers may augment initial risk stratification and individualization of patient care" "The widely used k-top scoring pair (k-TSP) algorithm is a simple yet powerful parameter-free classifier. It owes its success in many cancer microarray datasets to an effective feature selection algorithm that is based on relative expression ordering of gene pairs. However, its general robustness does not extend to some difficult datasets, such as those involving cancer outcome prediction, which may be due to the relatively low voting scheme used by the classifier. We believe that the performance can be enhanced by separating its effective feature selection component and combining it with a powerful classifier such as the support vector machine (SVM). [...] We developed an approach integrating the k-TSP ranking algorithm (TSP) with other machine learning methods, allowing combination of the computationally efficient, multivariate feature ranking of k-TSP with multivariate classifiers such as SVM. We evaluated this hybrid scheme (k-TSP+SVM) in a range of simulated datasets with known data structures. As compared with other feature selection methods, such as a univariate method similar to Fisher's discriminant criterion (Fisher) or a recursive feature elimination embedded in SVM (RFE), TSP is increasingly more effective than the other two methods as the informative genes become progressively more correlated, which is demonstrated both in terms of the classification performance and the ability to recover true informative genes." "Integrating genomic information with traditional clinical risk factors to improve the prediction of disease outcomes could profoundly change the practice of medicine. However, the large number of potential markers and possible complexity of the relationship between markers and disease make it difficult to construct accurate risk prediction models. [...] In recent years, much work has been done to group genes into pathways and networks. Integrating such biological knowledge into statistical learning could potentially improve model interpretability and reliability. One effective approach is to employ a kernel machine (KM) framework, which can capture nonlinear effects of nonlinear kernels are used. [...] In this article, we derive testing and prediction methods for KM regression under the accelerated failure time (AFT) model, a useful alternative to the PH model. We approximate the null distribution of our test statistic using resampling procedures. When multiple kernels are of potential interest, it may be unclear in advance which kernel to use for testing and estimation. We propose a robust Omnibus Test that combines information across kernels, and an approach for selecting the best kernel for estimation. The methods are illustrated with an application in breast cancer."
181	Shi, P and Ray, S and Zhu, Q and Kang, M A	Top scoring pairs for feature selection in machine learning and applications to cancer outcome prediction	Bmc Bioinform	12		15-15	2011	USA	"4 cancer prognosis microarray datasets, including data with > 50 samples per group	Case-control study	error rate (LOOCV, test set)	cross-validation + test set	"The widely used k-top scoring pair (k-TSP) algorithm is a simple yet powerful parameter-free classifier. It owes its success in many cancer microarray datasets to an effective feature selection algorithm that is based on relative expression ordering of gene pairs. However, its general robustness does not extend to some difficult datasets, such as those involving cancer outcome prediction, which may be due to the relatively low voting scheme used by the classifier. We believe that the performance can be enhanced by separating its effective feature selection component and combining it with a powerful classifier such as the support vector machine (SVM). [...] We developed an approach integrating the k-TSP ranking algorithm (TSP) with other machine learning methods, allowing combination of the computationally efficient, multivariate feature ranking of k-TSP with multivariate classifiers such as SVM. We evaluated this hybrid scheme (k-TSP+SVM) in a range of simulated datasets with known data structures. As compared with other feature selection methods, such as a univariate method similar to Fisher's discriminant criterion (Fisher) or a recursive feature elimination embedded in SVM (RFE), TSP is increasingly more effective than the other two methods as the informative genes become progressively more correlated, which is demonstrated both in terms of the classification performance and the ability to recover true informative genes." "Integrating genomic information with traditional clinical risk factors to improve the prediction of disease outcomes could profoundly change the practice of medicine. However, the large number of potential markers and possible complexity of the relationship between markers and disease make it difficult to construct accurate risk prediction models. [...] In recent years, much work has been done to group genes into pathways and networks. Integrating such biological knowledge into statistical learning could potentially improve model interpretability and reliability. One effective approach is to employ a kernel machine (KM) framework, which can capture nonlinear effects of nonlinear kernels are used. [...] In this article, we derive testing and prediction methods for KM regression under the accelerated failure time (AFT) model, a useful alternative to the PH model. We approximate the null distribution of our test statistic using resampling procedures. When multiple kernels are of potential interest, it may be unclear in advance which kernel to use for testing and estimation. We propose a robust Omnibus Test that combines information across kernels, and an approach for selecting the best kernel for estimation. The methods are illustrated with an application in breast cancer."
182	Simnett, J A and Cai, T	Omnibus risk assessment via accelerated failure time kernel machine modeling	Biometrics	69	4	861-873	2013	United States	"training set of 454 lymph node negative breast cancer patients [...] A total of 119 deaths or recurrences were observed"	Case-control study	C-statistic (training + validation data)	cross-validation + test set	"The parathyroid glands are located adjacent to the thyroid and occasionally within it. Enlarged and hyperplastic parathyroid glands can be mistaken as thyroid nodules or suspicious thyroid nodules. On fine needle aspiration biopsy (FNAB) of such lesions, cytology is often nondiagnostic, failing to identify its parathyroid origin and potentially resulting in an unnecessary thyroid surgery. The Affirma Genomic Sequencing Classifier (GSC) identifies genomically benign thyroid nodules among those with indeterminate FNAB to prevent unnecessary diagnostic surgery using RNA sequencing and machine learning algorithms. [...] The final classifier was blindly tested on an independent test set of 195 FNAs (118 Bethesda II, 77 Bethesda IV). The classifier had a sensitivity of 41% for thyroid correctly called positive (CI=8-100%) and 100% specificity (191/191 thyroid correctly called negative; CI=98.1-100%)." "For predicting relevant clinical outcomes, we propose a flexible statistical machine learning approach that acknowledges and models the interaction between platform-specific measurements through nonlinear kernel machines and borrows information within and between platforms through a hierarchical Bayesian framework. Our model has parameters with direct interpretations in terms of the effects of platforms and data interactions within and across platforms. The parameter estimation algorithm in our model uses a computationally efficient variational Bayes approach that scales well to large high-throughput datasets. [...] We apply our methods of integrating gene/mRNA expression and microRNA profiles for predicting patient survival times to the Cancer Genome Atlas (TCGA) based glioblastoma multiforme (GBM) dataset. In terms of prediction accuracy, we show that our non-linear and interaction-based integrative methods perform better than linear alternatives and non-integrative methods that do not account for interactions between the platforms. We also find several prognostic mRNAs and microRNAs that are related to tumor invasion and are known to drive tumor metastasis and severe inflammatory response in GBM. [...] Our approach gains its flexibility and power by modeling the non-linear interaction structures between and within the platforms." "Machine learning (ML) may harbor the potential to capture the metabolic complexity in Alzheimer Disease (AD). Here we set out to test the performance of metabolites in blood to categorize AD when compared to CSF biomarkers. [...] Deep Learning (DL), Extreme Gradient Boosting (XGBoost) and Random Forest (RF) were used to differentiate AD from CN. These models were internally validated using Nested Cross Validation (NCV). [...] On the test data, DL produced the AUC of 0.85 (0.80-0.89), XGBoost produced 0.88 (0.86-0.89) and RF produced 0.85 (0.83-0.87). By comparison, CSF measures of amyloid, p-tau and tau (together with age and gender) produced with XGBoost the AUC values of 0.78, 0.83 and 0.87, respectively. [...] This study showed that plasma metabolites have the potential to match the AUC of well-established AD CSF biomarkers in a relatively small cohort." " [...] We performed a systematic and comprehensive evaluation of several major algorithms for multicategory classification, several gene selection methods, multiple ensemble classifier methods and two cross-validation designs using 11 datasets spanning 74 diagnostic categories and 41 cancer types and 12 normal tissue types. [...] Multicategory support vector machines (MC-SVMs) are the most effective classifiers in performing accurate cancer diagnosis from gene expression data. The MC-SVM techniques by Crammer and Singer, Weston and Watkins and one-versus-rest were found to be the best methods in this domain. MC-SVMs outperform other popular machine learning algorithms, such as k-nearest neighbors, backpropagation and probabilistic neural networks, often to a remarkable degree. Gene selection techniques significantly improve the classification performance of both MC-SVMs and other non-SVM learning algorithms. Ensemble classifiers do not generally improve performance of the best non-ensemble models. These results guided the construction of a software system GEMS (Gene Expression Model Selector) that automates high-quality model construction and enforces sound optimization and performance evaluation procedures."
183	Whitney, D and Voh, M and Ladenson, P W	Clinical validation of the Affirma genomic sequencing parathyroid classifier	Thyroid	27		ASO-AS1	2017		"476 FNAs-6 parathyroid and 470 thyroid FNAs	Case-control study	sensitivity, specificity (training and validation cohort)	cross-validation + external cohort validation	"The parathyroid glands are located adjacent to the thyroid and occasionally within it. Enlarged and hyperplastic parathyroid glands can be mistaken as thyroid nodules or suspicious thyroid nodules. On fine needle aspiration biopsy (FNAB) of such lesions, cytology is often nondiagnostic, failing to identify its parathyroid origin and potentially resulting in an unnecessary thyroid surgery. The Affirma Genomic Sequencing Classifier (GSC) identifies genomically benign thyroid nodules among those with indeterminate FNAB to prevent unnecessary diagnostic surgery using RNA sequencing and machine learning algorithms. [...] The final classifier was blindly tested on an independent test set of 195 FNAs (118 Bethesda II, 77 Bethesda IV). The classifier had a sensitivity of 41% for thyroid correctly called positive (CI=8-100%) and 100% specificity (191/191 thyroid correctly called negative; CI=98.1-100%)." "For predicting relevant clinical outcomes, we propose a flexible statistical machine learning approach that acknowledges and models the interaction between platform-specific measurements through nonlinear kernel machines and borrows information within and between platforms through a hierarchical Bayesian framework. Our model has parameters with direct interpretations in terms of the effects of platforms and data interactions within and across platforms. The parameter estimation algorithm in our model uses a computationally efficient variational Bayes approach that scales well to large high-throughput datasets. [...] We apply our methods of integrating gene/mRNA expression and microRNA profiles for predicting patient survival times to the Cancer Genome Atlas (TCGA) based glioblastoma multiforme (GBM) dataset. In terms of prediction accuracy, we show that our non-linear and interaction-based integrative methods perform better than linear alternatives and non-integrative methods that do not account for interactions between the platforms. We also find several prognostic mRNAs and microRNAs that are related to tumor invasion and are known to drive tumor metastasis and severe inflammatory response in GBM. [...] Our approach gains its flexibility and power by modeling the non-linear interaction structures between and within the platforms." "Machine learning (ML) may harbor the potential to capture the metabolic complexity in Alzheimer Disease (AD). Here we set out to test the performance of metabolites in blood to categorize AD when compared to CSF biomarkers. [...] Deep Learning (DL), Extreme Gradient Boosting (XGBoost) and Random Forest (RF) were used to differentiate AD from CN. These models were internally validated using Nested Cross Validation (NCV). [...] On the test data, DL produced the AUC of 0.85 (0.80-0.89), XGBoost produced 0.88 (0.86-0.89) and RF produced 0.85 (0.83-0.87). By comparison, CSF measures of amyloid, p-tau and tau (together with age and gender) produced with XGBoost the AUC values of 0.78, 0.83 and 0.87, respectively. [...] This study showed that plasma metabolites have the potential to match the AUC of well-established AD CSF biomarkers in a relatively small cohort." " [...] We performed a systematic and comprehensive evaluation of several major algorithms for multicategory classification, several gene selection methods, multiple ensemble classifier methods and two cross-validation designs using 11 datasets spanning 74 diagnostic categories and 41 cancer types and 12 normal tissue types. [...] Multicategory support vector machines (MC-SVMs) are the most effective classifiers in performing accurate cancer diagnosis from gene expression data. The MC-SVM techniques by Crammer and Singer, Weston and Watkins and one-versus-rest were found to be the best methods in this domain. MC-SVMs outperform other popular machine learning algorithms, such as k-nearest neighbors, backpropagation and probabilistic neural networks, often to a remarkable degree. Gene selection techniques significantly improve the classification performance of both MC-SVMs and other non-SVM learning algorithms. Ensemble classifiers do not generally improve performance of the best non-ensemble models. These results guided the construction of a software system GEMS (Gene Expression Model Selector) that automates high-quality model construction and enforces sound optimization and performance evaluation procedures."
184	Srivastava, S and Wang, W and Manjani, G and Ordones, C and Baladandayyaiah, V	Integrating multi-platform genomic data using hierarchical Bayesian relevance vector machines	Biometrics	2013	1		2013		"GBM data have multiple molecular measurements on over 500 samples that include gene expression, copy number, methylation and microRNA expression"	Case-control study	mean square prediction error ("We randomly split the GBM survival data into a training data and a test data with 223 (90%) and 25 (10%) patients, respectively")	training + test set	" [...] We performed a systematic and comprehensive evaluation of several major algorithms for multicategory classification, several gene selection methods, multiple ensemble classifier methods and two cross-validation designs using 11 datasets spanning 74 diagnostic categories and 41 cancer types and 12 normal tissue types. [...] Multicategory support vector machines (MC-SVMs) are the most effective classifiers in performing accurate cancer diagnosis from gene expression data. The MC-SVM techniques by Crammer and Singer, Weston and Watkins and one-versus-rest were found to be the best methods in this domain. MC-SVMs outperform other popular machine learning algorithms, such as k-nearest neighbors, backpropagation and probabilistic neural networks, often to a remarkable degree. Gene selection techniques significantly improve the classification performance of both MC-SVMs and other non-SVM learning algorithms. Ensemble classifiers do not generally improve performance of the best non-ensemble models. These results guided the construction of a software system GEMS (Gene Expression Model Selector) that automates high-quality model construction and enforces sound optimization and performance evaluation procedures."
185	Legido-Quigley, C	A comprehensive evaluation of multicategory classification methods for microarray gene expression cancer diagnosis	Biometrics	21	5	631-643	2005	USA	"11 datasets spanning 74 diagnostic categories and 41 cancer types and 12 normal tissue types	Case-control study	accuracy, relative classifier information (Design 1: nested stratified 10-fold CV outer loop, 9-fold CV inner loop, Design 2: nested LOOCV outer loop, 10-fold CV inner loop)	cross-validation	" [...] We performed a systematic and comprehensive evaluation of several major algorithms for multicategory classification, several gene selection methods, multiple ensemble classifier methods and two cross-validation designs using 11 datasets spanning 74 diagnostic categories and 41 cancer types and 12 normal tissue types. [...] Multicategory support vector machines (MC-SVMs) are the most effective classifiers in performing accurate cancer diagnosis from gene expression data. The MC-SVM techniques by Crammer and Singer, Weston and Watkins and one-versus-rest were found to be the best methods in this domain. MC-SVMs outperform other popular machine learning algorithms, such as k-nearest neighbors, backpropagation and probabilistic neural networks, often to a remarkable degree. Gene selection techniques significantly improve the classification performance of both MC-SVMs and other non-SVM learning algorithms. Ensemble classifiers do not generally improve performance of the best non-ensemble models. These results guided the construction of a software system GEMS (Gene Expression Model Selector) that automates high-quality model construction and enforces sound optimization and performance evaluation procedures."

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46

<p>Wang, Y and Cao, Y and Song, W and He, K and Li, J and Wang, J and Xu, B and Shi, H and Yu, H and C and 209 Li, A, L</p>	<p>Serum peptide pattern that differentially diagnoses hepatitis B virus-related hepatocellular carcinoma from liver cirrhosis</p> <p>Gastroenterology and Hepatology</p>	<p>29</p>	<p>7</p>	<p>1544-1550</p>	<p>2014</p>	<p>China</p>	<p>http://dx.doi.org/10.1111/jgh.12648</p>	<p>article</p>	<p>80 HCC and 67 LC patients</p> <p>Case-control study</p> <p>AUC, accuracy, sensitivity, specificity (10-fold cross-validation + test set)</p> <p>cross-validation + test set</p>	<p>"Although alpha-fetoprotein (AFP) is a useful serologic marker of hepatocellular carcinoma (HCC), it is not sufficiently sensitive to differentiate HCC and liver cirrhosis (LC) caused by hepatitis B virus (HBV) infection. [...] With a highly optimized peptide extraction and matrix-assisted laser desorption/ionization time of flight/time of flight mass spectrometry approach, we investigated serum peptide profiles of 80 HCC and 67 LC patients. Three supervised machine learning methods were employed to construct classifiers [...] We proposed a novel method for distinguishing HCC from cirrhosis, based on a multilayer perceptron (MLP) method. We obtained a sensitivity of 90.0%, specificity of 79.4%, and overall accuracy of 85.1% on an independent test set. The combination of the MLP method and serum AFP level performed serum AFP marker alone in distinguishing HCC patients from LC patients."</p>	<p>"Although alpha-fetoprotein (AFP) is a useful serologic marker of hepatocellular carcinoma (HCC), it is not sufficiently sensitive to differentiate HCC and liver cirrhosis (LC) caused by hepatitis B virus (HBV) infection. [...] With a highly optimized peptide extraction and matrix-assisted laser desorption/ionization time of flight/time of flight mass spectrometry approach, we investigated serum peptide profiles of 80 HCC and 67 LC patients. Three supervised machine learning methods were employed to construct classifiers [...] We proposed a novel method for distinguishing HCC from cirrhosis, based on a multilayer perceptron (MLP) method. We obtained a sensitivity of 90.0%, specificity of 79.4%, and overall accuracy of 85.1% on an independent test set. The combination of the MLP method and serum AFP level performed serum AFP marker alone in distinguishing HCC patients from LC patients."</p>
<p>Wang, S and Li, M, C</p>	<p>Impacts of Predictive Genomic Classifier Performance on Subpopulation-Specific Treatment Effects Assessment</p> <p>Statistics in Biosciences</p>	<p>8</p>	<p>1</p>	<p>129-158</p>	<p>2016</p>	<p></p>	<p>http://dx.doi.org/10.1007/s12561-016-9103-6</p>	<p>article</p>	<p>80 HCC and 67 LC patients</p> <p>Case-control study</p> <p>accuracy, sensitivity, specificity, PPV, NPV, permutation p-values (cross-validation + external validation)</p> <p>cross-validation + external cohort validation</p>	<p>"We investigate the classification performance characteristics of a binary genomic composite biomarker (expected to be predictive of treatment effects) including sensitivity, specificity, accuracy, positive predictive value and negative predictive value as a function of true sensitive prevalence. In doing so, we report the finding based on three representative tuning parameter sets with varying degree of rigor in their choices of the parameters ranging from highly rigorous, moderately rigorous to mildly rigorous, to articulate the rationales on the choices of tuning parameter sets. We also study the impacts of misclassification of genomic biomarker classifiers on their assessment of treatment effects in the positive and negative patient subpopulations, and all corner patients."</p>	<p>"We investigate the classification performance characteristics of a binary genomic composite biomarker (expected to be predictive of treatment effects) including sensitivity, specificity, accuracy, positive predictive value and negative predictive value as a function of true sensitive prevalence. In doing so, we report the finding based on three representative tuning parameter sets with varying degree of rigor in their choices of the parameters ranging from highly rigorous, moderately rigorous to mildly rigorous. We articulate the rationales on the choices of tuning parameter sets. We also study the impacts of misclassification of genomic biomarker classifiers on their assessment of treatment effects in the positive and negative patient subpopulations, and all corner patients."</p>
<p>Wei, Z and Wang, K and Ou, H and Zhang, H and Bradford, J and Kim, C and Frackelton, E and Hou, C and Glessner, T and Chikawa, R and Stanley, C and Mones, D and Grant, S F and Polychronakos, C and Hakonarson, H</p>	<p>From disease association to risk assessment: an optimistic view from genome-wide association studies on type 1 diabetes</p> <p>PLoS Genetics</p>	<p>5</p>	<p>10</p>	<p>e1000678</p>	<p>2009</p>	<p>United States of America</p>	<p>http://dx.doi.org/10.1371/journal.pgen.1000678</p>	<p>article</p>	<p>80 HCC and 67 LC patients</p> <p>Case-control study</p> <p>AUC, accuracy, sensitivity, specificity (5-fold cross-validation + test set)</p> <p>cross-validation + test set</p>	<p>"We accessed the 500K Affymetrix chip genotype data from WTCCC on ~1,500 samples from the 1958 British Birth Cohort, ~1,500 samples from the UK Blood Service Control Group, as well as ~2,000 samples each from the following disease collections: type 1 diabetes (T1D), type 2 diabetes (T2D), rheumatoid arthritis (RA), inflammatory bowel disease (IBD), bipolar disorder (BD), hypertension (HT), coronary artery disease (CAD)"</p>	<p>"We accessed the 500K Affymetrix chip genotype data from WTCCC on ~1,500 samples from the 1958 British Birth Cohort, ~1,500 samples from the UK Blood Service Control Group, as well as ~2,000 samples each from the following disease collections: type 1 diabetes (T1D), type 2 diabetes (T2D), rheumatoid arthritis (RA), inflammatory bowel disease (IBD), bipolar disorder (BD), hypertension (HT), coronary artery disease (CAD)"</p>
<p>White, B S and Khan, S A and Ammad-Ud-Din, M and Potdar, S and Mason, M J and Tognon, C E and Draker, B J and Hechtma, C A and Kulkarni, O P and Kurtz, S E and Parkka, K and Tyner, J W and 212 Attikolaki, T and Wernberg, K and Gunney, J</p>	<p>Gene expression predicts ex vivo drug sensitivity in acute myeloid leukemia</p> <p>Cancer Research</p>	<p>78</p>	<p>13</p>	<p></p>	<p>2018</p>	<p></p>	<p>http://dx.doi.org/10.1158/1538-7443.2018.01333</p>	<p>meeting abstract</p>	<p>80 HCC and 67 LC patients</p> <p>Case-control study</p> <p>correlation, p-value (10-fold CV)</p> <p>cross-validation</p>	<p>"We harmonized two large-scale AML ex vivo studies screened for drug response and profiled transcriptionally—OHSU (803 AML patient samples and 140 drugs) and FIMM (48 AML samples and 480 drugs)"</p>	<p>"We harmonized two large-scale AML ex vivo studies screened for drug response and profiled transcriptionally—OHSU (803 AML patient samples and 140 drugs) and FIMM (48 AML samples and 480 drugs)"</p>
<p>Wu, H and Cai, L and Li, D and Wang, X and Zhao, S and Zou, F and Zhou, K</p>	<p>Metagenomics Biomarkers Selected for Prediction of Three Different Diseases in Chinese Population</p> <p>Biomed Res Int</p>	<p>2018</p>	<p></p>	<p>2936257</p>	<p>2018</p>	<p>China</p>	<p>http://dx.doi.org/10.1155/2018/2936257</p>	<p>article</p>	<p>80 HCC and 67 LC patients</p> <p>Case-control study</p> <p>AUC, F1 score (5-fold CV)</p> <p>cross-validation</p>	<p>"microbiome of 806 Chinese individuals (383 controls, 170 with type 2 diabetes, 130 with rheumatoid arthritis, and 123 with liver cirrhosis)"</p>	<p>"microbiome of 806 Chinese individuals (383 controls, 170 with type 2 diabetes, 130 with rheumatoid arthritis, and 123 with liver cirrhosis)"</p>
<p>Xie, G and Dong, C and Kong, Y and Guoqing, J F and Li, M and Wang, K</p>	<p>Group lasso regularized deep learning for cancer prognosis from multi-omics and clinical features</p> <p>Genes</p>	<p>10</p>	<p>3</p>	<p></p>	<p>2019</p>	<p>USA</p>	<p>http://dx.doi.org/10.3390/genes10030240</p>	<p>article</p>	<p>80 HCC and 67 LC patients</p> <p>Case-control study</p> <p>the used dataset cover more than 50 samples per group</p> <p>cross-validation</p>	<p>"To harness the rich information in multi-omics data, we developed GDP (Group lasso regularized deep learning for cancer prognosis), a computational tool for survival prediction using both clinical and multi-omics data. GDP integrated a deep learning framework and Cox proportional hazard model (CPH) together, and applied group lasso regularization to incorporate gene-level group prior knowledge into the model training process. We evaluated its performance in both simulated and real data from The Cancer Genome Atlas (TCGA) project."</p>	<p>"To harness the rich information in multi-omics data, we developed GDP (Group lasso regularized deep learning for cancer prognosis), a computational tool for survival prediction using both clinical and multi-omics data. GDP integrated a deep learning framework and Cox proportional hazard model (CPH) together, and applied group lasso regularization to incorporate gene-level group prior knowledge into the model training process. We evaluated its performance in both simulated and real data from The Cancer Genome Atlas (TCGA) project."</p>
<p>Yang, J S and Zhu, Z and He, S and Ji, Z and Hee</p>	<p>Minimal-redundancy-maximal-relevance feature selection using different relevance measures for omics data classification</p> <p>Computational Intelligence in Bioinformatics and Computational Biology</p>	<p>2013</p>	<p>251</p>	<p></p>	<p>2013</p>	<p>USA</p>	<p>http://dx.doi.org/10.1109/CIBCB.2013.6560417</p>	<p>article</p>	<p>80 HCC and 67 LC patients</p> <p>Case-control study</p> <p>accuracy (10 runs of 10-fold CV)</p> <p>cross-validation</p>	<p>"This paper studies a minimal-redundancy-maximal-relevance (MRMR) feature selection for omics data classification using three different relevance evaluation measures including mutual information (MI), correlation coefficient (CC), and maximal information coefficient (MIC). A linear forward search method is used to search the optimal feature subset. The experimental results on five real-world omics datasets indicate that MRMR feature selection with CC is more robust to obtain better (or competitive) classification accuracy than the other two measures."</p>	<p>"This paper studies a minimal-redundancy-maximal-relevance (MRMR) feature selection for omics data classification using three different relevance evaluation measures including mutual information (MI), correlation coefficient (CC), and maximal information coefficient (MIC). A linear forward search method is used to search the optimal feature subset. The experimental results on five real-world omics datasets indicate that MRMR feature selection with CC is more robust to obtain better (or competitive) classification accuracy than the other two measures."</p>
<p>Yang, S and Naiman, D Q</p>	<p>Multiclass cancer classification based on gene expression comparison</p> <p>Stat Appl Genet Mol Biol</p>	<p>13</p>	<p>4</p>	<p>477-496</p>	<p>2014</p>	<p>United States</p>	<p>http://www.jstor.org/stable/24775737</p>	<p>article</p>	<p>80 HCC and 67 LC patients</p> <p>Case-control study</p> <p>accuracy (LOOCV + test set)</p> <p>cross-validation + test set</p>	<p>"one of the considered datasets has more than 50 samples per group (ML) is a two-stage study where a retrospective stage 1 generated expression profiles for 2,143 patients and was designed for biomarker discovery. A prospective stage 2 produced an independent cohort of 1,152 patients and was used for validation"</p>	<p>"one of the considered datasets has more than 50 samples per group (ML) is a two-stage study where a retrospective stage 1 generated expression profiles for 2,143 patients and was designed for biomarker discovery. A prospective stage 2 produced an independent cohort of 1,152 patients and was used for validation"</p>
<p>Yang, Y and Huang, N and Hao, L and Kong, W</p>	<p>A clustering-based approach for efficient identification of microRNA combinatorial biomarkers</p> <p>BMC Genomics</p>	<p>18</p>	<p></p>	<p>210-210</p>	<p>2017</p>	<p>China</p>	<p>http://dx.doi.org/10.1186/s12864-017-3458-8</p>	<p>article</p>	<p>80 HCC and 67 LC patients</p> <p>Case-control study</p> <p>accuracy, sensitivity, specificity (5-fold cross-validation)</p> <p>cross-validation</p>	<p>"This study aims to select combinatorial miRNA biomarkers, which have higher sensitivity and specificity than single-gene biomarkers. In order to avoid exhaustive search and redundant information, miRNAs are firstly clustered, then the combinations of representative cluster members are assessed as potential biomarkers. [...] Our experimental results demonstrate that the clustering-based method can identify microRNA combinatorial biomarkers with high accuracy and efficiency"</p>	<p>"This study aims to select combinatorial miRNA biomarkers, which have higher sensitivity and specificity than single-gene biomarkers. In order to avoid exhaustive search and redundant information, miRNAs are firstly clustered, then the combinations of representative cluster members are assessed as potential biomarkers. [...] Our experimental results demonstrate that the clustering-based method can identify microRNA combinatorial biomarkers with high accuracy and efficiency"</p>

218	Yu, H and Samuels, D and Zhao, Y and Guo, Y	Architectures and accuracy of artificial neural networks for disease classification from omics data	BMC Genomics	20		12-12	2019	China	http://dx.doi.org/10.1186/s12854-019-0566-z	article	37 omics datasets, including datasets with more than 50 samples per group	Case-control study	Accuracy, Cohen's Kappa (nested cross-validation)	cross-validation	<p>"Deep learning has made tremendous successes in numerous artificial intelligence applications and is unsurprisingly penetrating into various biomedical domains. High-throughput omics data in the form of molecular profile matrices, such as transcriptional and metabolomic data, have long existed as a valuable resource for facilitating diagnosis of patient statuses/stages. It is timely imperative to compare deep learning neural networks against classical machine learning methods in the setting of matrix-formatted omics data in terms of classification accuracy and robustness. [...] Using 37 high-throughput omics datasets, covering transcriptomes and metabolomes, we evaluated the classification power of deep learning compared to traditional machine learning methods. Representative deep learning methods, Multi-Layer Perceptrons (MLP) and Convolutional Neural Networks (CNN), were deployed and explored in seeking optimal architectures for the best classification performance. Together with five classical supervised classification methods (Linear Discriminant Analysis, Multinomial Logistic Regression, Naive Bayes, Random Forest, Support Vector Machine), MLP and CNN were comparatively tested on the 37 datasets to predict disease stages or to discriminate diseased samples from normal samples. MLP achieved the highest overall accuracy among all methods tested. More thorough analyses revealed that single hidden layer MLPs with ample hidden units outperformed deeper MLPs. Furthermore, MLP was one of the most robust methods against imbalanced class composition and inaccurate class labels. [...] Our results concluded that shallow MLPs (of one or two hidden layers) with ample hidden neurons are sufficient to achieve superior and robust classification performance in exploring numerical matrix-formatted omics data for diagnosis purpose."</p> <p>"Non-small cell lung cancer (NSCLC) is the most common type of lung cancer while adenocarcinoma (LUAD) is its most common subtype. [...] However, it remains difficult to find the most significant genetic features and build a high-effective predictive model for treatment outcomes. To confront the challenges, we collected the large-scale LUAD case data with both genome and clinic information (n = 371) from TCGA (The Cancer Genome Atlas) (http://cancergenome.nih.gov), analyzed the somatic mutation difference between the two groups categorized based on the 3-year overall survival, and developed a machine learning model to predict prognosis based on the most significant genetic markers. Through the analysis, we identified a list of genes with different mutation frequencies between different prognosis groups and many were involved in cell-cell adhesion and motility; an absolute majority of the genes showed higher mutation frequencies in the good prognosis group."</p> <p>"Ovarian cancer is the deadliest gynecologic malignancy in the United States, with most patients diagnosed in the advanced stage of the disease. Platinum-based antineoplastic therapeutics is indispensable to treating advanced ovarian serous carcinoma. However, patients have heterogeneous response to platinum drugs, and it is difficult to predict these inter-individual differences before administering medication. In this study, we investigated the tumor proteomic profiles and clinical characteristics of 130 ovarian serous carcinoma patients analyzed by the Clinical Proteomic Tumor Analysis Consortium (CPTAC), predicted the platinum drug response using supervised machine learning methods [...]. Our data-driven feature selection approach indicated that tumor proteomic profiles contain information for predicting platinum response (P<0.0001). We further built a least absolute shrinkage and selection operator (LASSO) Cox proportional hazards model that stratified patients into early relapse and late relapse groups (P=0.00013). The top proteomic features indicative of platinum response were involved in ATP synthesis pathways and Ras GTPase binding."</p> <p>"Prostate cancer (PCa) is the second leading cause of cancer-related mortality in men. The prevalent diagnosis method is based on the serum Prostate-Specific Antigen (PSA) screening test, which suffers from low specificity, over-diagnosis and overtreatment. In this work, untargated metabolomic profiling of age-matched serum samples from prostate cancer patients and healthy individuals was performed using ultra-performance liquid chromatography coupled to high-resolution tandem mass spectrometry (UPLC-MS/MS) and machine learning methods. A metabolite-based in vitro diagnostic multivariate index assay (VDMA) was developed to predict the presence of PCa in serum samples with high classification sensitivity, specificity and accuracy. A panel of 40 metabolic spectral features was found to be differential with 92.1% sensitivity, 94.3% specificity, and AUC=0.9703 accuracy. The performance of the VDMA was higher than the prevalent PSA test. [...] The identification of fatty acids, amino acids, liposphospholipids, and bile acids provided further insights into the metabolic alterations associated with the disease."</p> <p>"High-risk neuroblastoma is a very aggressive disease, with excessive tumor growth and poor outcome. A proper stratification of the high-risk patients by prognostic outcome is important for treatment. However, there is still a lack of survival stratification for the high-risk neuroblastoma. To fill the gap, we adopt a deep learning algorithm, Autoencoder, to integrate multi-omics data, and combine it with K-means clustering to identify two subtypes with significant survival differences. By comparing the Autoencoder with PCA, Cluster, and tDScore about the classification based on multi-omics data integration, Autoencoder-based classification outperforms the alternative approaches. Furthermore, we also validated the classification in two independent datasets by training machine-learning classification models, and confirmed its robustness. Functional analysis revealed that MYCN amplification was more frequently associated with the high-risk subtype in accordance with the overexpression of MYCN/MYC in this subtype."</p> <p>"Prostate cancer is a leading male malignancy worldwide, while the prognosis prediction remains quite inaccurate. The study aimed to observe whether there was an association between the prognosis of prostate cancer and genetic mutation profile, and to build an accurate prognostic predictor based on the genetic signatures. [...] No significant gene with somatic mutation rate difference was found between prognostic groups of prostate cancer. Total 43 atypical genes were screened for building a support vector machine model to predict prostate cancer prognosis, with an average accuracy of 68% and 64% for 5-fold cross-validation or training-testing evaluation respectively. When combined with the National Institute for Health and Care Excellence (NICE) features, the model could be further improved, with the 5-fold cross-validation accuracy of 71%, much better than NICE itself (62%)."</p> <p>"Primary platinum-based chemoresistance occurs in approximately one-third of patients with serous ovarian cancer (SOC); however, traditional clinical indicators are poor predictors of chemoresistance. So we aimed to identify novel genes as predictors of primary platinum-based chemoresistance. Gene expression microarray analyses were performed to identify the genes related to primary platinum resistance in SOC on two discovery datasets (GSE51373, GSE63885) and one validation dataset (TCGA). Univariate and multivariate analysis with logistic regression were performed to evaluate the predictive values of the genes for platinum resistance. Machine learning algorithms (linear kernel support vector machine and artificial neural network) were applied to build predictive models. Univariate and multivariate analyses with Cox proportional hazards regression and log-rank tests were used to assess the effect of these genetic signatures for platinum resistance on prognosis in two independent datasets (GSE51373, GSE63885). AGGF1 and MFR4 were found to be highly expressed in patients with platinum-resistant SOC and independently predicted platinum resistance. Platinum resistance prediction models based on these targets had robust predictive power (Highest AUC: 0.8056, 95% CI: 0.6388-0.9773; lowest AUC: 0.7245, 95% CI: 0.6388-0.9773).</p> <p>"Despite existing prognostic markers, breast cancer prognosis remains a difficult subject due to the complex relationships between many contributing factors and survival. This study seeks to integrate multiple clinicopathological and genomic factors with dimensional reduction across machine learning algorithms to compare survival predictions. [...] ROC and accuracy were not significantly different between models (ROC and accuracy around 0.67 and 0.72 across models, respectively). However, ensemble methods resulted in better fit (CS) with stable measures of variable importance across 10 random training/validation splits. K-means clustering of gene expression profiles on training data points along with KNN classification of validation data points was a robust method of dimensional reduction. Furthermore, the gene cluster with the highest mortality risk was an influential factor in model prediction [...]. Using machine learning methods to construct predictive models for 5-year survival in patients with breast cancer, we demonstrated discrimination ability across models with new insight into the stability and utility of dimensional reduction on genomic features in breast cancer survival prediction. In this review, we reviewed the most recent published works that used deep learning to build models for cancer prognosis prediction. Deep learning has been suggested to be a more generic model, requires less data engineering, and achieves more accurate prediction when working with large amounts of data. The application of deep learning in cancer prognosis has been shown to be equivalent or better than current approaches, such as Cox-Ph. With the burst of multi-omics data, including genomic data, transcriptomic data and clinical information in cancer studies, we believe that deep learning would potentially improve cancer prognosis."</p>
219	Yu, J and Hu, Y and Xu, Y and Wang, J and Kuang, J and Zhang, W and Shao, J and Guo, D and Wang, Y	LUADs: an effective prediction model on prognosis of lung adenocarcinoma based on somatic mutational features	BMC Cancer	19	1	263-263	2019	China	http://dx.doi.org/10.1186/s12929-019-0431-7	article	more than 50 samples per group for "good" vs. "poor" distinction (Table S3)	Case-control study	ROC, accuracy, sensitivity, specificity (5-fold CV)	cross-validation	<p>"The TARGE1 cohort is comprised of 407 high-risk neuroblastoma samples, including 217 samples with gene expression data and 380 samples with copy number alterations (CNA). Among these obtained samples, 180 had both gene expression and CNA data. The SEQ cohort has a total of 498 neuroblastoma samples, including 176 high-risk and 322 low- or intermediate-risk samples."</p>
220	Zhou, Y, K H and Levine, D A and Zhang, H and Chan, D W and Zhang, Z and Snyder, M	Predicting Ovarian Cancer Patients' Clinical Response to Platinum-Based Chemotherapy by Their Tumor Proteomic Signatures	J Proteome Res	15	8	2455-2465	2018	United States	http://dx.doi.org/10.1021/acs.jproteome.8b00119	article	"Proteomic profiles of 130 ovarian serous carcinoma patients were analyzed by The Cancer Genome Atlas (TCGA) Clinical Proteomic Tumor Analysis Consortium (CPTAC)."	Cases only (drug response study)	AUC (LOOCV, hold-out test set)	training + test set	<p>"Prostate cancer (PCa) is the second leading cause of cancer-related mortality in men. The prevalent diagnosis method is based on the serum Prostate-Specific Antigen (PSA) screening test, which suffers from low specificity, over-diagnosis and overtreatment. In this work, untargated metabolomic profiling of age-matched serum samples from prostate cancer patients and healthy individuals was performed using ultra-performance liquid chromatography coupled to high-resolution tandem mass spectrometry (UPLC-MS/MS) and machine learning methods. A metabolite-based in vitro diagnostic multivariate index assay (VDMA) was developed to predict the presence of PCa in serum samples with high classification sensitivity, specificity and accuracy. A panel of 40 metabolic spectral features was found to be differential with 92.1% sensitivity, 94.3% specificity, and AUC=0.9703 accuracy. The performance of the VDMA was higher than the prevalent PSA test. [...] The identification of fatty acids, amino acids, liposphospholipids, and bile acids provided further insights into the metabolic alterations associated with the disease."</p> <p>"High-risk neuroblastoma is a very aggressive disease, with excessive tumor growth and poor outcome. A proper stratification of the high-risk patients by prognostic outcome is important for treatment. However, there is still a lack of survival stratification for the high-risk neuroblastoma. To fill the gap, we adopt a deep learning algorithm, Autoencoder, to integrate multi-omics data, and combine it with K-means clustering to identify two subtypes with significant survival differences. By comparing the Autoencoder with PCA, Cluster, and tDScore about the classification based on multi-omics data integration, Autoencoder-based classification outperforms the alternative approaches. Furthermore, we also validated the classification in two independent datasets by training machine-learning classification models, and confirmed its robustness. Functional analysis revealed that MYCN amplification was more frequently associated with the high-risk subtype in accordance with the overexpression of MYCN/MYC in this subtype."</p>
221	Zang, X and Jones, C M and Long, T Q and Monge, M E and Zhou, M and Walker, L D and Mezeniec, R and Gray, A and McDonald, J F and Fernandez, F M	Feasibility of detecting prostate cancer by ultraperformance liquid chromatography-mass spectrometry serum metabolomics	J Proteome Res	13	7	3444-3454	2014	USA	http://dx.doi.org/10.1021/acs.jproteome.4b0004a	article	Age-matched blood serum samples were obtained from 64 PCA patients (age range 49-65, mean age 59 ± 4 year) and 50 healthy individuals (age range 45-76, mean age 57.7 years).	Case-control study	accuracy, sensitivity, specificity (10-fold cross-validation)	cross-validation	<p>"Prostate cancer is a leading male malignancy worldwide, while the prognosis prediction remains quite inaccurate. The study aimed to observe whether there was an association between the prognosis of prostate cancer and genetic mutation profile, and to build an accurate prognostic predictor based on the genetic signatures. [...] No significant gene with somatic mutation rate difference was found between prognostic groups of prostate cancer. Total 43 atypical genes were screened for building a support vector machine model to predict prostate cancer prognosis, with an average accuracy of 68% and 64% for 5-fold cross-validation or training-testing evaluation respectively. When combined with the National Institute for Health and Care Excellence (NICE) features, the model could be further improved, with the 5-fold cross-validation accuracy of 71%, much better than NICE itself (62%)."</p>
222	Zhang, L and Lv, C and Jin, Y and Cheng, G and Fu, Y and Yuan, D and Tao, Y and Guo, Y and Ni, X and Shi, T	Deep learning based multi-omics data analysis for predicting prognosis of high-risk neuroblastoma	Frontiers in Genetics	9			2018	China	http://dx.doi.org/10.3389/fgene.2018.00111	article	The TARGE1 cohort has a total of 498 neuroblastoma samples, including 176 high-risk and 322 low- or intermediate-risk samples.	Cases only (prognostic stratification)	C-index, log rank p-value, AUC (10-fold CV, external validation set)	cross-validation + external cohort validation	<p>"Prostate cancer is a leading male malignancy worldwide, while the prognosis prediction remains quite inaccurate. The study aimed to observe whether there was an association between the prognosis of prostate cancer and genetic mutation profile, and to build an accurate prognostic predictor based on the genetic signatures. [...] No significant gene with somatic mutation rate difference was found between prognostic groups of prostate cancer. Total 43 atypical genes were screened for building a support vector machine model to predict prostate cancer prognosis, with an average accuracy of 68% and 64% for 5-fold cross-validation or training-testing evaluation respectively. When combined with the National Institute for Health and Care Excellence (NICE) features, the model could be further improved, with the 5-fold cross-validation accuracy of 71%, much better than NICE itself (62%)."</p> <p>"Primary platinum-based chemoresistance occurs in approximately one-third of patients with serous ovarian cancer (SOC); however, traditional clinical indicators are poor predictors of chemoresistance. So we aimed to identify novel genes as predictors of primary platinum-based chemoresistance. Gene expression microarray analyses were performed to identify the genes related to primary platinum resistance in SOC on two discovery datasets (GSE51373, GSE63885) and one validation dataset (TCGA). Univariate and multivariate analysis with logistic regression were performed to evaluate the predictive values of the genes for platinum resistance. Machine learning algorithms (linear kernel support vector machine and artificial neural network) were applied to build predictive models. Univariate and multivariate analyses with Cox proportional hazards regression and log-rank tests were used to assess the effect of these genetic signatures for platinum resistance on prognosis in two independent datasets (GSE51373, GSE63885). AGGF1 and MFR4 were found to be highly expressed in patients with platinum-resistant SOC and independently predicted platinum resistance. Platinum resistance prediction models based on these targets had robust predictive power (Highest AUC: 0.8056, 95% CI: 0.6388-0.9773; lowest AUC: 0.7245, 95% CI: 0.6388-0.9773).</p> <p>"Despite existing prognostic markers, breast cancer prognosis remains a difficult subject due to the complex relationships between many contributing factors and survival. This study seeks to integrate multiple clinicopathological and genomic factors with dimensional reduction across machine learning algorithms to compare survival predictions. [...] ROC and accuracy were not significantly different between models (ROC and accuracy around 0.67 and 0.72 across models, respectively). However, ensemble methods resulted in better fit (CS) with stable measures of variable importance across 10 random training/validation splits. K-means clustering of gene expression profiles on training data points along with KNN classification of validation data points was a robust method of dimensional reduction. Furthermore, the gene cluster with the highest mortality risk was an influential factor in model prediction [...]. Using machine learning methods to construct predictive models for 5-year survival in patients with breast cancer, we demonstrated discrimination ability across models with new insight into the stability and utility of dimensional reduction on genomic features in breast cancer survival prediction. In this review, we reviewed the most recent published works that used deep learning to build models for cancer prognosis prediction. Deep learning has been suggested to be a more generic model, requires less data engineering, and achieves more accurate prediction when working with large amounts of data. The application of deep learning in cancer prognosis has been shown to be equivalent or better than current approaches, such as Cox-Ph. With the burst of multi-omics data, including genomic data, transcriptomic data and clinical information in cancer studies, we believe that deep learning would potentially improve cancer prognosis."</p>
223	Zhang, S and Xu, Y and Hul, X and Yang, F and Hu, Y and Shao, J and Liang, H and Wang, Y	Improvement in prediction of prostate cancer prognosis with somatic mutational signatures	Journal of Cancer	8	16	3261-3267	2017	China	http://dx.doi.org/10.7196/jco.1291	article	more than 50 samples per group (both for recurrence status and tumor status)	Cases only (prognosis study)	ROC, accuracy (5-fold CV, training/test split)	training + test set	<p>"Prostate cancer is a leading male malignancy worldwide, while the prognosis prediction remains quite inaccurate. The study aimed to observe whether there was an association between the prognosis of prostate cancer and genetic mutation profile, and to build an accurate prognostic predictor based on the genetic signatures. [...] No significant gene with somatic mutation rate difference was found between prognostic groups of prostate cancer. Total 43 atypical genes were screened for building a support vector machine model to predict prostate cancer prognosis, with an average accuracy of 68% and 64% for 5-fold cross-validation or training-testing evaluation respectively. When combined with the National Institute for Health and Care Excellence (NICE) features, the model could be further improved, with the 5-fold cross-validation accuracy of 71%, much better than NICE itself (62%)."</p>
224	Zhao, H and Sun, Q and Li, L and Zhou, J and Zhang, C and Hu, T and Zhou, X and Zhang, L and Wang, B and Li, L and Shi, T and Li, H	High expression levels of AGGF1 and MFR4 are associated with adverse prognosis in patients with serous ovarian cancer	Journal of Cancer	10	2	397-407	2019	China	http://dx.doi.org/10.7196/jco.18172	article	The used TCGA data covers more than 50 samples per group	Cases only (drug resistance prediction)	log rank test p-value, AUC (5-times 10-fold CV + external validation)	cross-validation + external cohort validation	<p>"Despite existing prognostic markers, breast cancer prognosis remains a difficult subject due to the complex relationships between many contributing factors and survival. This study seeks to integrate multiple clinicopathological and genomic factors with dimensional reduction across machine learning algorithms to compare survival predictions. [...] ROC and accuracy were not significantly different between models (ROC and accuracy around 0.67 and 0.72 across models, respectively). However, ensemble methods resulted in better fit (CS) with stable measures of variable importance across 10 random training/validation splits. K-means clustering of gene expression profiles on training data points along with KNN classification of validation data points was a robust method of dimensional reduction. Furthermore, the gene cluster with the highest mortality risk was an influential factor in model prediction [...]. Using machine learning methods to construct predictive models for 5-year survival in patients with breast cancer, we demonstrated discrimination ability across models with new insight into the stability and utility of dimensional reduction on genomic features in breast cancer survival prediction. In this review, we reviewed the most recent published works that used deep learning to build models for cancer prognosis prediction. Deep learning has been suggested to be a more generic model, requires less data engineering, and achieves more accurate prediction when working with large amounts of data. The application of deep learning in cancer prognosis has been shown to be equivalent or better than current approaches, such as Cox-Ph. With the burst of multi-omics data, including genomic data, transcriptomic data and clinical information in cancer studies, we believe that deep learning would potentially improve cancer prognosis."</p>
225	Zhao, M and Tang, Y and Kim, H and Hasagawa, K	Machine Learning With K-Means Dimensional Reduction for Predicting Survival Outcomes in Patients With Breast Cancer	Cancer Informatics	17			2018	USA	http://dx.doi.org/10.1177/1533033818770116	article	"2509 adult female participants with breast cancer in a prospective cohort study"	Cases only (survival prediction)	ROC, accuracy (10-fold CV)	cross-validation	<p>"Despite existing prognostic markers, breast cancer prognosis remains a difficult subject due to the complex relationships between many contributing factors and survival. This study seeks to integrate multiple clinicopathological and genomic factors with dimensional reduction across machine learning algorithms to compare survival predictions. [...] ROC and accuracy were not significantly different between models (ROC and accuracy around 0.67 and 0.72 across models, respectively). However, ensemble methods resulted in better fit (CS) with stable measures of variable importance across 10 random training/validation splits. K-means clustering of gene expression profiles on training data points along with KNN classification of validation data points was a robust method of dimensional reduction. Furthermore, the gene cluster with the highest mortality risk was an influential factor in model prediction [...]. Using machine learning methods to construct predictive models for 5-year survival in patients with breast cancer, we demonstrated discrimination ability across models with new insight into the stability and utility of dimensional reduction on genomic features in breast cancer survival prediction. In this review, we reviewed the most recent published works that used deep learning to build models for cancer prognosis prediction. Deep learning has been suggested to be a more generic model, requires less data engineering, and achieves more accurate prediction when working with large amounts of data. The application of deep learning in cancer prognosis has been shown to be equivalent or better than current approaches, such as Cox-Ph. With the burst of multi-omics data, including genomic data, transcriptomic data and clinical information in cancer studies, we believe that deep learning would potentially improve cancer prognosis."</p>
226	Zhou, W and Xie, L and Han, J and Guo, X	The application of deep learning in cancer prognosis prediction	Cancers	12	3				http://dx.doi.org/10.3390/cancers12030354	review					https://doi.org/10.3390/cancers12030354

1	Bergamaschi, A and Ku, J and Ning, Y and Collin, F and Ellison, C and Phillips, T and McCarthy, I and Wang, W and Antino, M and Haan, D and Scott, A and Loyd, P and Puler, G and Ashworth, A and 245 Quake, S and Lev, S	Epicomic detection of multiple cancers in plasma derived cell free DNA	Cancer Research	80	16	2020 USA	meeting abstract	11810.1136/bmjopen-2020-025178	Case-control study	AUC (5-fold CV)	cross-validation	"These findings suggest that ShmC changes in cfDNA enable non-invasive detection of early stage breast, pancreatic, prostate, and lung cancers. Furthermore, ShmC profiling of cfDNA may enable the prediction of clinically relevant features such as tumor size in early-stage adenocarcinoma or indolent disease in prostate cancer. Finally, this study identifies a suite of ShmC biomarkers that may be further validated in larger, and more diverse, patient cohorts."
2	Berry, S E and Valdes, A M and Drew, D A and Aniczar, F and Maziou, M and Wolf, J and Capdevilla, J and Hallgrangsson, G and Davies, R and Al Khalib, H and Bonnett, C and Ganeva, S and Bakker, E and Hart, D and Mangino, M and Merino, J and Linberg, L and Wyatt, P and Oudovas, M and Gardner, C	Human postprandial responses to food and potential for precision nutrition	Nat Med	26	6	964-973	2020 UK	10.1038/s41591-020-0814-1	Response to food intake prediction	correlation	training + test set	"We developed a machine learning model that predicted both triglyceride ($r = 0.47$) and glycemic ($r = 0.77$) responses to food intake. These findings may be informative for developing personalized diet strategies."
3	246 D and Dehalayth, L M and Chan, A T and Segata, N and Franks, P W and Spector, T D	Exposition based biomarkers and models to classify early and late-stage samples of Papillary Thyroid Carcinoma	PLoS One	15	4	e0231629	2020 India	10.1371/journal.pone.0231629	Case-control study	AUC (cross-validation + external validation)	cross-validation + test set	"We developed a machine learning model that predicted both triglyceride ($r = 0.47$) and glycemic ($r = 0.77$) responses to food intake. These findings may be informative for developing personalized diet strategies."
4	247 Bhillal, S and Kaur, H and Kaur, S and Sharma, S and Raghava, G P S	"paired tumor and normal whole-exome sequencing (WES) data from clinically annotated tumor specimens treated with anti-CTLA4 (n=14), melanoma) and anti-PD1 therapies (n=94, KCLC melanoma, head and neck, bladder cancer, and others)"	Cancer Research	80	16	2020 USA	meeting abstract	11810.1136/bmjopen-2020-025178	Case-control study	precision, recall ("A random forest classifier was trained and tested on various subsets of the dataset")	training + test set	"We developed a machine learning model that predicted both triglyceride ($r = 0.47$) and glycemic ($r = 0.77$) responses to food intake. These findings may be informative for developing personalized diet strategies."
5	248 Bigelow, E G and Baraz, A and Yarchuan, M and Jaffee, E M	Machine learning methods to identify clinical genomic predictors of clinical responses to immune checkpoint inhibitor therapy	Cancer Research	80	16	2020 USA	meeting abstract	11810.1136/bmjopen-2020-025178	Case-control study	precision, recall ("A random forest classifier was trained and tested on various subsets of the dataset")	training + test set	"We developed a machine learning model that predicted both triglyceride ($r = 0.47$) and glycemic ($r = 0.77$) responses to food intake. These findings may be informative for developing personalized diet strategies."
6	Brown, E and Karrar, A and Hellings, S and Stepanova, M and Warrack, B and Lam, B and Onorato, J and Felix, S and Aghajani, A and Jeffers, T and Rajput, A and Charles, E and Nader, F and Luo, Y and Rely, 249 M and Zhao, L and Thompson, J and Goodman, Z and Youmans, Z	Metabonomics composite biomarkers selected by machine learning predicts identification of a Tumor Microenvironment relevant Gene Set based on gastric cancer gene expression datasets from 1659 patients from five independent cohorts"	Theranostics	5	10	19 8633-8647	2020 China	10.1039/c9th24733a	Prognostic study	concordance index (C-index; was calculated to determine the discrimination of the nomogram via a bootstrap method with 1000 resamples)	external cohort validation	"As a tumor microenvironment-relevant gene set based prognostic signature, the GP5C model provides an effective approach to evaluate GC [Gastric Cancer] patient survival outcomes, and may prolong overall survival by enabling the selection of individualized targeted therapy."
7	250 Cai, W Y and Dong, Z N and Fu, X T and Liu, L Y and Wang, L and Ye, G D and Luo, Q C and Chen, Y C	Gut microbiome, big data and machine learning to promote precision medicine for cancer	Nat Rev Gastroenterol Hepatol	17	10	635-648	2020 Italy	10.1038/s41575-020-0092-3	Review	review (not applicable)	review	"In this Perspective, we provide a brief overview on the role of gut microbiome in cancer and focus on the need, role and limitations of a machine learning-driven approach to analyze large amounts of complex health-care information in the era of big data."
8	Cammarota, G and Ianiro, G and Ahera, A and Carbone, C and Temko, A and Claesson, M J and 251 Gasbarrini, A and Tortora, G	"A 230-sample training set was extracted randomly from the TCGA dataset, representing the same 60/40 ER/ER- proportion of the PAM50 test set. All the remaining 597 cases were instead included in the TCGA test set."	Sci Rep	10	1	14071	2020 Italy	10.1038/s41598-020-70312-2	Tumor subtype prediction	AUC (10-fold CV, training/test set)	cross-validation + external cohort validation	"This, in conclusion, the main contribution of this paper is twofold: 1. Propose the AWCFA reference reconstruction approach to face the proved issues of the standard PAM50 algorithm. 2. Define RNA-seq based classification approaches to perform single-sample BC intrinsic subtyping with external AWCFA-based PAM50 or regularized mlR methods. These strategies appeared valuable to favor the use of RNA-seq in BC [Breast Cancer] clinical practice and are worthy of other studies on heterogeneous RNA-seq data, to evaluate and strengthen the reliability of their intrinsic subtyping methods."
9	252 Cascianelli, S and Molteni, I and Isella, C and Maseroli, M and Medico, E	Machine learning for RNA sequencing-based intrinsic subtyping of breast cancer	Sci Rep	10	1	14071	2020 Italy	10.1038/s41598-020-70312-2	Tumor subtype prediction	AUC (10-fold CV, training/test set)	cross-validation + external cohort validation	"Although standard therapy affected every gene signature and significantly increased myeloid cell signatures, logistic regression analysis determined that ancestral background significantly changed 23 of 34 gene signatures. Additionally, the strongest association to gene expression changes was found with autoantibodies, and this also had etiology in ancestry: the AA predisposition to have both HbF and d/dNA autoantibodies compared with AA predisposition to have only anti-d/dNA. A machine learning approach was used to determine a gene signature characteristic to distinguish AA-SL and was most influenced by genes characteristic of the perturbed B cell axis in AA-SL patients."
10	253 Cascianelli, S and Molteni, I and Isella, C and Maseroli, M and Medico, E	Machine learning for RNA sequencing-based intrinsic subtyping of breast cancer	Sci Rep	10	1	14071	2020 Italy	10.1038/s41598-020-70312-2	Tumor subtype prediction	AUC (10-fold CV, training/test set)	cross-validation + external cohort validation	"Although standard therapy affected every gene signature and significantly increased myeloid cell signatures, logistic regression analysis determined that ancestral background significantly changed 23 of 34 gene signatures. Additionally, the strongest association to gene expression changes was found with autoantibodies, and this also had etiology in ancestry: the AA predisposition to have both HbF and d/dNA autoantibodies compared with AA predisposition to have only anti-d/dNA. A machine learning approach was used to determine a gene signature characteristic to distinguish AA-SL and was most influenced by genes characteristic of the perturbed B cell axis in AA-SL patients."
11	Catalina, M D and Bachali, P and Yeo, A E and Geraci, N S and Petri, M A and Grammer, A C and 254 Lipky, E	Patient ancestry significantly contributes to molecular heterogeneity of systemic lupus erythematosus	JCI Insight	5	15	2020 USA	meeting abstract	11810.1136/bmjopen-2020-025178	Subtype categorization	AUC (10-fold CV)	cross-validation	"After correction for repeated measures, clustering identified 3 separate metabolic/clinical profiles: 1) right ventricular (RV) dysfunction, arrhythmia and degree (n=107), 2) complex biventricular disease with hypoxia and lower educational level (n=79) and 3) individuals managing well with valvular and septal defects (n=42). Metabonomic data permitted the creation of models associated with prevalent arrhythmia (cross-validated AUC 0.67), patient reported exertion-associated dyspnea of breath (AUC 0.58) and RV dysfunction (AUC 0.62)."
12	Cedars, A M and Manthorpe, C and Ko, J and Bottiglieri, T and Weingarten, A and Opostowsky, A and 254 Kutys, S	ARTIFICIAL INTELLIGENCE FACILITATED METABOLOMIC PROFILING IN ADULT CONGENITAL HEART DISEASE (ACHD)	Journal of the American College of Cardiology	75	11	552-562	2020 USA	10.1016/j.jacc.2020.03.026	Subtype categorization	AUC (cross-validation)	cross-validation	"While ML methods are powerful, they are still similar to previous clinical tools in that physician interpretation is crucial for implementation in a real-world setting. We should be cognizant of how potential biases can interfere with the black box nature of these algorithms. It is also important to make these technologies inclusive of skin of color. Further research in ML should be transparent by making algorithms and datasets available to the public for further validation and testing."
13	255 Chan, S and Reddy, V and Myers, B and Thibodeaux, Q and Brownstone, N and Liao, W	Machine Learning in Dermatology: Current Applications, Opportunities, and Limitations	Dermatol Ther	10	3	365-386	2020 USA	10.1016/j.jaad.2020.03.026	Review	review (not applicable)	review	"Urology is a constantly changing specialty with a wide range of therapeutic breakthroughs. A huge understanding of the genomic expression profiles for each urological cancer and a tendency to use cutting-edge technology to treat our patients. All of the major developments must be analyzed objectively, taking into account costs to the health systems, risks and benefits to the patients, and the legal background that comes with these. A critical analysis of these new technologies and pharmacological breakthroughs should be made before considering changing our clinical practice."
14	256 Chavarraga, J and Moreno, C	Precision Medicine, Artificial Intelligence, and Genetic Markers in Urology: Do we need to Tailor our Clinical Practice?	Urologia Colombiana	29	3	158-167	2020 Colombia	10.1016/j.uro.2020.03.026	Review	review (not applicable)	review	"Although it is clear that no single method is consistently preferable, and that most of the proposed approaches are task and/or data dependent, the complexity of tumor analysis suggests that network-based approaches are needed. In this context, it is clear that omics integration is one of the most promising and demanding challenges of the modern informatics, and that there is an urgent need to prove the reproducibility, interpretability, and generalization capability of the proposed methods."
15	Cherici, M and Bussola, N and Marcolini, A and Francescatti, M and Zandonà, A and Trastulla, L and 257 Agostinelli, C and Jurman, G and Furlanello, C	Integrative Network Fusion: A Multi-Omics Approach in Molecular Profiling	Frontiers in Oncology	10	2020 Italy	meeting abstract	11810.1136/bmjopen-2020-025178	Disease status, subtype and survival prediction	Matthews Correlation Coefficient (10 x stratified Monte Carlo cross-validation (50% training/validation proportion))	cross-validation + external cohort validation	"The immune mechanism underlying acute food-allergic events remains elusive until today. Deciphering this immunological response shall enable to identify novel biomarker for stratification of patients into reaction endotypes. The availability of powerful multi-omics technologies, together with integrated data analysis, network based approaches and unbiased machine learning holds out the prospect of providing clinically useful biomarkers or biomarker signatures being predictive for reaction phenotypes."	
16	Czolk, R and Klueber, J and Sarsenen, M and Wilmes, P and Codreanu Morel, F and Skov, P S and 258 Hillger, C and Bindlev-Jensen, C and Ollert, M and Kuehn, A	lglf-Mediated Peanut Allergy: Current and Novel Predictive Biomarkers for Clinical Phenotypes Using Multi-Omics Front	Immunol Front	11	594350-594350	2020 g	Luxembourg	10.3389/fimmu.2020.00460	Review	review (not applicable)	review	"Paracoccidiodiomycosis (PCM) is a fungal infection typically found in Latin American countries, especially in Brazil. The identification of this disease is based on techniques that may fall sometimes. Intending to improve PCM detection in patient samples, this study used the combination of two of the newest technologies, artificial intelligence and metabolomics. This combination allowed PCM detection, independently of disease form, through identification of a set of molecules present in patient's blood. The great difference in this research was the ability to detect disease with better confidence than the routine methods employed today."
17	de Oliveira Lima, E and Navarro, L C and Morishita, K N and Kamikawa, C M and Rodrigues, R M and Dabaja, M Z and Ribeiro, D M and DeLencastre, J and Dias-Audebert, F L and de Silva Ribeiro, M and 259 Vicentini, A P and Rocha, A and Catharino, R R	Metabonomics and machine learning approaches combined in pursuing more accurate paracoccidiodiomycosis diagnoses	mSystems	5	3	2020 Brazil	meeting abstract	11810.1136/bmjopen-2020-025178	Case-control study	AUC (cross-validation)	cross-validation	"Technological advances now enable large-scale datasets, including DNA and RNA sequence data, proteomics and metabolomics data, to be captured from individuals and groups of patients along the genotype-phenotype continuum of chronic kidney disease (CKD). The ability to combine these high-dimensional datasets, in which the number of variables exceeds the number of clinical outcome observations, using computational approaches such as machine learning, provides an opportunity to re-classify patients into molecularly defined subgroups that better reflect underlying disease mechanisms. Patients with CKD are uniquely poised to benefit from these integrative, multi-omics approaches since the kidney, urinary, blood and urine samples used to generate these different types of molecular data are frequently obtained during routine clinical care."
18	260 Eddy, S and Mariani, L H and Kretzel, M	Integrated multi-omics approaches to improve classification of chronic kidney disease	Nat Rev Nephrol	16	11	657-668	2020 USA	10.1038/s41581-020-0296-5	Review	review (not applicable)	review	"Newly proposed biomarkers offer precise and noninvasive ways to monitor patient's status. Cell-free DNA quantitation is increasingly explored as an indicator of allograft injury and rejection, which can help avoid unnecessary biopsies and more frequently monitor graft function."
19	261 Fu, S and Zarringer, A	Recent advances in precision medicine for individualized immunosuppression	Curr Opin Organ Transplant	25	4	420-425	2020 USA	10.1097/COO.0000000000000771	Review	review (not applicable)	review	"We show that features of gut microbiome, in combination with already used fecal biomarkers, are strong predictors for differentiating IBD and IBS, with additional potential of classifying location type of IBD. These results have a potential to improve non-invasive pre-screening for IBD in clinical practice."
20	262 Kurlishov, A and Fu, J and Zhenkova, I and Weersma, R	Microbiome and fecal biomarkers can diagnose and classify inflammatory bowel disease	Netherlands Journal of Gastroenterology	7	8	166-167	2019 s	10.1016/j.njg.2019.04.001	Case-control study	AUC, sensitivity, specificity (independent validation)	training + test set	"We show that features of gut microbiome, in combination with already used fecal biomarkers, are strong predictors for differentiating IBD and IBS, with additional potential of classifying location type of IBD. These results have a potential to improve non-invasive pre-screening for IBD in clinical practice."
21	263 Garcia, S and Lauritsen, J and Zhang, C and Dagaard, M O and Nielsen, R L and Daugaard, G and 263 Gupta, R	Prediction of nephrotoxicity associated with cisplatin-based chemotherapy in testicular cancer patients	JNCI Cancer Spectrum	4	3	2020 USA	meeting abstract	11810.1136/bmjopen-2020-025178	Case-control study	AUC (Training and testing of the algorithm was performed with a 5 group cross-validation)	cross-validation + external cohort validation	"We show that features of gut microbiome, in combination with already used fecal biomarkers, are strong predictors for differentiating IBD and IBS, with additional potential of classifying location type of IBD. These results have a potential to improve non-invasive pre-screening for IBD in clinical practice."

1	Gindin, Y and Chang, J and Billi, A and Camargo, M and Huss, R and Chung, C and Myers, P and P and 264 Younossi, S and Harrison, S and Anstee, Q and M and Lombard, R	Hepatology	72	1	43A-44A	2020	USA	https://doi.org/10.1007/s12250-020-00187-1	meeting abstract	"The study included 1,120 adults with advanced fibrosis (F3-F4) due to NAFLD enrolled in the interventional and STELAR trials (discovery cohort, n=94) and the ATLAS trial (validation cohort, n=126)..."	Case-control study	AUC (training + validation cohort)	training + test validation	"A machine learning technique applied to hepatic transcriptomic data identified a 30-gene expression signature that accurately differentiates NAFLD patients with cirrhosis from those with bridging fibrosis. The functional activities of these genes may suggest novel drivers of fibrosis progression in NAFLD."
2										"We obtained 273 TEP expression profiles spanning six cancer types: non-small-cell lung cancer (NSCLC); 59, colorectal cancer (CRC); 44, glioblastoma multiforme (GBM); 40, breast cancer (BRCA); 38, pancreatic cancer (PC); 33, hepatocellular carcinoma (HCC). In addition to the cancer samples, platelets from 54 healthy individuals were also profiled."	Case-control study	AUC (LOOCV)	cross-validation	"In this article, we demonstrated the predictive power of a small set of platelet genes in determining the existence of cancer. Similar strategies can be developed for inferring the potential cancer types. In all these cases, the gene panels need to be validated on larger patient and control samples' cohorts. An orthogonal application of such panels could be tracking the treatment responses, as well as the recurrence of the disease."
3	Goswami, C and Chawla, S and Thakral, D and Pant, H and Verma, P and Malik, P and Jayadeva and 265 Gupta, S and Anshu, G and Sengupta, D	Genomics	21	1	744-744	2020	India	https://doi.org/10.1007/s12250-020-00187-1	article	Molecular signature comprising 11 platelet genes enables accurate blood-BMC based diagnosis of NSCLC	Case-control study	AUC (LOOCV)	cross-validation	"Artificial intelligence has great potential to advance diagnosis and treatment of patients with neurocognitive disorders. Multi-feature datasets can improve personalization and predictive ability of machine learning algorithms in healthcare. Development of Explainable Artificial Intelligence is warranted to establish trust in models for clinical decision-making."
4										"The experiment of this study is conducted on a public dataset of microarray prostate cancer gene expression, consisting of 102 tissue samples (52 prostate tumor and 50 normal tissues)".	Case-control study	accuracy (10-fold CV)	cross-validation	"Artificial intelligence has great potential to advance diagnosis and treatment of patients with neurocognitive disorders. Multi-feature datasets can improve personalization and predictive ability of machine learning algorithms in healthcare. Development of Explainable Artificial Intelligence is warranted to establish trust in models for clinical decision-making."
5	Graham, S and Lee, E E and Jeste, D V and Van Patten, R and Twamley, E W and Nebeker, C and 266 Yamada, Y and Kim, H C and Deppe, C A	Psychiatry Research	284			2020	USA	https://doi.org/10.1016/j.psychres.2020.151322	article	Artificial intelligence approaches to predicting and detecting cognitive decline in older adults: A conceptual review	Review	Review	Review	"In this paper, we propose to use a correlation feature selection (CFS) method with random committee (RC) ensemble learning to detect prostate cancer from microarray data of gene expression. A set of experiments are conducted on a public benchmark dataset using 10-fold cross-validation technique to evaluate the proposed approach. The experimental results revealed that the proposed approach attains 95.08% accuracy, which is higher than related work methods on the same dataset."
6										"The experiment of this study is conducted on a public dataset of microarray prostate cancer gene expression, consisting of 102 tissue samples (52 prostate tumor and 50 normal tissues)".	Case-control study	accuracy (10-fold CV)	cross-validation	"In this paper, we propose to use a correlation feature selection (CFS) method with random committee (RC) ensemble learning to detect prostate cancer from microarray data of gene expression. A set of experiments are conducted on a public benchmark dataset using 10-fold cross-validation technique to evaluate the proposed approach. The experimental results revealed that the proposed approach attains 95.08% accuracy, which is higher than related work methods on the same dataset."
7	267 Gumeal, A and Sam Samudoo, R and Al-Rakhami, M and AlSalami, H and El-Zaart, A	Health Informatics Journal	27	1		2021	Arabia	https://doi.org/10.1177/146342522095221	article	Feature selection with ensemble learning for prostate cancer diagnosis from microarray gene expression	Case-control study	accuracy (10-fold CV)	cross-validation	"We identified 34 biomarkers and 19 KEGG pathways associated with ovarian cancer. The independent test results in three GEO datasets proved the robustness of our model. The literature reviewing show 19 (56%) biomarkers and 8 (42.1%) KEGG pathways identified based on the classification subtypes have been proved to be associated with ovarian cancer."
8										"The R package TCGA-assemble2 [13] was used for data collection and we obtained 798 samples concluded three types of omics data: mRNA-seq data (LNC, Illumina HiSeq, RNASeq V2), copy number data (CCSC, Illumina HiSeq) and copy number variation (CNV) data (BROAD-MIT Genome wide SNP_6)".	Tumor stratification	silhouette score (external test datasets)	training + test set	"Precision medicine matches each individual to the best treatment in a way that is tailored to his/her genetic uniqueness. To further personalize medicine, precision medicine increasingly incorporates and integrates data beyond genomics, such as epigenomics and metabolomics, as well as imaging. Increasingly, the robust use and integration of these modalities in precision medicine require the use of artificial intelligence and machine learning. This modern view of precision medicine, adopted early in certain areas of medicine such as cancer, has started to impact the field of reproductive medicine."
9	268 Guo, L Y and Wu, A H and Wang, Y X and Zhang, L P and Chi, H and Liang, X F	BioData Mining	13	1		2020	China	https://doi.org/10.1186/s13040-020-00272-x	article	Deep learning-based ovarian cancer subtypes identification using multi-omics data	Case-control study	accuracy (10-fold CV)	cross-validation	"We integrated somatic mutations and previously used data types, including Exp, CNV, Methy, and protein, using MML to predict breast cancer patient survival. Applying mMR-selected features and MML classification, we found that the integration of somatic mutations enriched the diversity of features and was conducive to the improvement of the prediction model. In all, integrating promising data sources such as somatic mutations and harnessing the powerful feature selection method mMR and the effective data fusion method MML can increase the prediction accuracy of breast cancer patient survival."
10										"This study included a total of 202 subjects, comprising 82 AMD patients and 120 control subjects."	Case-control study	AUC (4-fold CV)	cross-validation	"Our study demonstrates that DNA methylation markers can robustly differentiate thyroid nodules based on their malignancy. They are thus promising candidates to develop non-invasive diagnostics for thyroid cancer screening"
11	269 Hajirasouliha, I and Simentoni, O	Fertil Steril	114	5	908-913	2020	USA	https://doi.org/10.1016/j.fertnstert.2020.09.016	article	Precision medicine and artificial intelligence: overview and relevance to reproductive medicine	Review	Review	Review	"Our study demonstrates that DNA methylation markers can robustly differentiate thyroid nodules based on their malignancy. They are thus promising candidates to develop non-invasive diagnostics for thyroid cancer screening"
12										"We finally obtained 488 primary breast tumors together with survival time, and all samples of them included all of the five aforementioned genomic data types. The details of our dataset are illustrated in Table 3."	Case-control study	AUC (entire datasets were randomly divided into a learning dataset (80% of the entire dataset) and validation dataset (20%))	training + test set	"Our study demonstrates that DNA methylation markers can robustly differentiate thyroid nodules based on their malignancy. They are thus promising candidates to develop non-invasive diagnostics for thyroid cancer screening"
13	Hao, S and Bai, J and Liu, H and Wang, L and Liu, T and Lin, C and Luo, X and Gao, J and Zhao, J and Liu, 270 H and Tang, H	Regenerative Therapy	15		180-186	2020	China	https://doi.org/10.1515/rtreg.2020.01501	article	Comparison of machine learning tools for the prediction of APO based on genetic, age, and diabetes-related variables in the Chinese population	Case-control study	AUC (4-fold CV)	cross-validation	"Our study demonstrates that DNA methylation markers can robustly differentiate thyroid nodules based on their malignancy. They are thus promising candidates to develop non-invasive diagnostics for thyroid cancer screening"
14										"We integrated somatic mutations and previously used data types, including Exp, CNV, Methy, and protein, using MML to predict breast cancer patient survival. Applying mMR-selected features and MML classification, we found that the integration of somatic mutations enriched the diversity of features and was conducive to the improvement of the prediction model. In all, integrating promising data sources such as somatic mutations and harnessing the powerful feature selection method mMR and the effective data fusion method MML can increase the prediction accuracy of breast cancer patient survival."	Case-control study	accuracy, sensitivity, specificity (training/test set)	training + test set+cohort	"Our study demonstrates that DNA methylation markers can robustly differentiate thyroid nodules based on their malignancy. They are thus promising candidates to develop non-invasive diagnostics for thyroid cancer screening"
15	271 He, Z and Zhang, J and Yuan, X and Zhang, Y	Frontiers in Genetics	11			2020	China	https://doi.org/10.3389/fgen.2020.01924	article	Integrating Somatic Mutations for Breast Cancer Survival Prediction Using Machine Learning Methods	Case-control study	AUC (entire datasets were randomly divided into a learning dataset (80% of the entire dataset) and validation dataset (20%))	training + test set+cohort	"Our study demonstrates that DNA methylation markers can robustly differentiate thyroid nodules based on their malignancy. They are thus promising candidates to develop non-invasive diagnostics for thyroid cancer screening"
16	Hong, S and Su, Z and Li, J and Yu, S and Lin, B and Ke, Z and Zhang, Q and Guo, Z and Lu, W and Peng, 272 S and Cheng, L and He, Q and Liu, R and Xiao, H	Annals of Oncology	31		51362-51362	2020	China	https://doi.org/10.1093/annonc/ndaa131	meeting abstract	Development of circulating free DNA methylation markers for thyroid nodule diagnostics	Case-control study	accuracy, sensitivity, specificity (training/test set)	training + test set+cohort	"Our study demonstrates that DNA methylation markers can robustly differentiate thyroid nodules based on their malignancy. They are thus promising candidates to develop non-invasive diagnostics for thyroid cancer screening"
17	Hoshino, A and Kim, H S and Bojmar, L and Gyan, K E and Coffi, M and Hernandez, J and Zambirini, C P and Rodrigues, G and Molina, H and Hestler, S and Mark, M T and Steiner, L and Santos-Martin, A and Lucetti, S and Di Giannantonio, A and Offer, K and Nakajima, M and Williams, C and Nogueis, L and Pelkser-Vatter, F and Hashimoto, S and Davies, E E and Freitas, D and Kenflick, C M and Arason, Y and Buehler, W and Lauritzen, P and Ogtuz, Y and Sugiura, K and Takahashi, N and Abdolvik, M and Bailey, K A and Jolissaint, J S and Wang, H and Harris, A and Schaeffer, L M and Garcia-Santos, G and Posner, Z and Batachandran, V P and Khakoo, Y and Raju, G P and Scherz, A and Sagi, I and Scherz-Shoval, R and Yarden, Y and Oren, M and Malabi, M and Petroncino, M and De Braganca, K C and Donzelli, M and Fischer, C and Vitelles, S and Wirth, G P and Gansbar, L and Marzano, M and Ahmed, A and DeStefano, J and Danzer, E and Roehrl, M H A and Lacayo, N J and Vincent, T C and Weiser, M R and Brady, M S and Meyers, P A and Weiser, L H and Ambati, S R and Chou, A J and Stoklin, E K and Modak, S and Roberts, S S and Bassi, E M and Dhillani, D and Krantz, B A and Cardoso, F and Simpson, A L and Berger, M and Rudin, C M and Simeone, D M and Jain, M and Ghajar, C M and Batra, S K and Stanger, B Z and Bai, J and Brown, K A and Rajasekhar, V K and Hralaj, J H and de Sousa, H and Kramer, K and Sheth, S and Banich, J and Pasqual, V and Heaton, T E and de Guellis, M and Pispas, D J and Schwarz, R and Zhang, H and Liu, Y and Shukla, A and Blawie, L and DeClerk, A and LaBarge, M and Bissell, M J and Caffrey, T C and Grandgett, P M and Hollingsworth, M A and Bromberg, J and Costa-Silva, B and Penhale, H and Kang, Y and Garcia, B A and O'Reilly, T M and 273 Keiser, D and Tripicci, T M and Jones, D R and Mawle, H and Jarrigan, W B and Lyden, D	Cell	182	4	1044-1044	2020	Japan	https://doi.org/10.1016/j.cell.2020.07.026	article	Extracellular Vesicles and Particle Biomarkers Define Multiple Human Cancers	Case-control study	sensitivity, specificity (10-fold CV + external test set)	cross-validation + external cohort validation	"Machine learning classification of plasma-derived EVP (extracellular vesicles and particles) cargo, including immunoglobulins, revealed 80% sensitivity/70% specificity in detecting cancer. Finally, we defined a panel of tumor-type-specific EVP proteins in TE5 and plasma, which can classify tumors of unknown primary origin. Thus, EVP proteins can serve as reliable biomarkers for cancer detection and determining cancer type."
18										"To confirm that EVPs are ideal diagnostic tools, we analyzed proteomes of TE5 (n=151) and plasma-derived (n=120) EVP."	Case-control study	sensitivity, specificity (10-fold CV + external test set)	cross-validation + external cohort validation	"Machine learning classification of plasma-derived EVP (extracellular vesicles and particles) cargo, including immunoglobulins, revealed 80% sensitivity/70% specificity in detecting cancer. Finally, we defined a panel of tumor-type-specific EVP proteins in TE5 and plasma, which can classify tumors of unknown primary origin. Thus, EVP proteins can serve as reliable biomarkers for cancer detection and determining cancer type."
19										"We investigated the two following of (4) involved 3,080 individuals (aged 32-81 years) examined between 2006 and 2008. For the second follow-up (FF4), 2,169 participants were examined from 2013 to 2014."	Case-control study	AUC (three-step feature selection with 100 random repeats of 10-fold cross validation)	cross-validation	"Our study demonstrates that DNA methylation markers can robustly differentiate thyroid nodules based on their malignancy. They are thus promising candidates to develop non-invasive diagnostics for thyroid cancer screening"
20	Huang, J and Huth, C and Govic, M and Troll, M and Adam, J and Zukoff, S and Prehn, C and Wang, L and Nano, J and Scherer, M F and Neshchen, S and Kastenmuller, G and Suhre, K and Layu, M and 274 Schlies, F and Gieger, C and Adamski, J and Hrabce de Angelis, M and Peters, A and Wang-Sattler, R	Diabetes	69		12 2756-2765	2020	Germany	https://doi.org/10.2337/db20-0566	article	Machine Learning Approaches Reveal Metabolic Signatures of Incident Chronic Kidney Disease in Individuals With Prediabetes and Type 2 Diabetes	Case-control study	AUC (three-step feature selection with 100 random repeats of 10-fold cross validation)	cross-validation	"Our study demonstrates that DNA methylation markers can robustly differentiate thyroid nodules based on their malignancy. They are thus promising candidates to develop non-invasive diagnostics for thyroid cancer screening"
21										"This performance comparison was conducted at pan-cancer level using 12 cancer from The Cancer Genome Atlas (TCGA). These 12 cancers were chosen due to their relatively large sample sizes and sufficient information about patient outcomes. The specific cancers analyzed in this paper were (1) Urothelial Bladder Carcinoma (BLCA); (2) Breast Invasive Carcinoma (BRCA); (3) Cervical Squamous Cell Carcinoma and Endocervical Adenocarcinoma (CESC); (4) Head Neck Squamous Cell Carcinoma (HNSC); (5) Kidney Renal Clear Cell Carcinoma (KIRC); (6) Kidney Renal Papillary Cell Carcinoma (KIPAP); (7) Liver Hepatocellular Carcinoma (LIHC); (8) Lung Adenocarcinoma (LUAD); (9) Lung Squamous Cell Carcinoma (LUSC); (10) Ovarian Cancer (OV); (11) Pancreatic Adenocarcinoma (PAAD); and (12) Stomach Adenocarcinoma (STAD). In this paper, we used the expression data of Illumina HiSeq RNA-seq V2 RSEM normalized genes from TCGA."	Cancer survival prognosis	C-index, p-value of log-rank test (Each dataset was split into training, validation, and testing sets in a proportion of 60, 20, and 20% respectively)	training + test set	"Overall our study demonstrated that the Deep Learning architecture can be effectively applied for cancer prognosis prediction with Cox-proportional hazard model incorporated. We found that Deep Learning based model demonstrated superior performances comparing to traditional machine learning models. Among the three Deep Learning based models tested, we observed that Co-net, which has the most succinct neural network structure, resulted in better prognosis performances in the measurement of concordance index and p-value of log-rank test. We showed that integrating autoencoder with Cox regression network does not significantly improve the prognosis performances."
22										"Serum of 520 subjects with HBVACLF (n=141), acute-on-chronic hepatic dysfunction (ACHD n=102), liver cirrhosis (LC n=110), chronic hepatitis B (CHB n=102), and normal controls (NC n=65) from a prospective multi-center cohort were subjected to a robust and highly streamlined single-run quantification proteomic analysis"	Cancer progression and prognosis prediction	AUC (training + validation cohort)	external cohort validation	"Overall our study demonstrated that the Deep Learning architecture can be effectively applied for cancer prognosis prediction with Cox-proportional hazard model incorporated. We found that Deep Learning based model demonstrated superior performances comparing to traditional machine learning models. Among the three Deep Learning based models tested, we observed that Co-net, which has the most succinct neural network structure, resulted in better prognosis performances in the measurement of concordance index and p-value of log-rank test. We showed that integrating autoencoder with Cox regression network does not significantly improve the prognosis performances."
23	Huang, J and Johnson, T S and Han, J and Helm, B and Cao, J and Zhang, C and Salama, P and Rizkalla, 275 M and Yu, C Y and Cheng, J and Xiang, S and Zhan, X and Zhang, J and Huang, C	BMC Med Genomics	13		41-41	2020	USA	https://doi.org/10.1186/s12920-020-00411-1	article	Deep learning-based cancer survival prognosis from RNA-seq data: approaches and evaluations	Cancer survival prognosis	C-index, p-value of log-rank test (Each dataset was split into training, validation, and testing sets in a proportion of 60, 20, and 20% respectively)	training + test set	"Overall our study demonstrated that the Deep Learning architecture can be effectively applied for cancer prognosis prediction with Cox-proportional hazard model incorporated. We found that Deep Learning based model demonstrated superior performances comparing to traditional machine learning models. Among the three Deep Learning based models tested, we observed that Co-net, which has the most succinct neural network structure, resulted in better prognosis performances in the measurement of concordance index and p-value of log-rank test. We showed that integrating autoencoder with Cox regression network does not significantly improve the prognosis performances."
24										"Quantitative proteomic study reveals an overall path of HBVACLF [Hepatitis B Virus-Related Acute-on-Chronic Liver Failure] disease progression. And the combinatorial predicted models provide fundamental information of multiple biomarkers response the disease progression, severity and prognosis."	Cancer progression and prognosis prediction	AUC (training + validation cohort)	external cohort validation	"Overall our study demonstrated that the Deep Learning architecture can be effectively applied for cancer prognosis prediction with Cox-proportional hazard model incorporated. We found that Deep Learning based model demonstrated superior performances comparing to traditional machine learning models. Among the three Deep Learning based models tested, we observed that Co-net, which has the most succinct neural network structure, resulted in better prognosis performances in the measurement of concordance index and p-value of log-rank test. We showed that integrating autoencoder with Cox regression network does not significantly improve the prognosis performances."
25	Jiang, J and Yao, H and Yang, L and Li, J and Xin, J and Shi, D and Liang, X and Cai, Q and Ren, K and 276 Chen, X and Li, J	Hepatology	70	suppl. 1	162A-163A	2019	China	https://doi.org/10.1093/ahp/afz046	meeting abstract	Proteome predicts progression and prognosis of hepatitis B virus-related acute-on-chronic liver failure	Cancer progression and prognosis prediction	AUC (training + validation cohort)	external cohort validation	"Overall our study demonstrated that the Deep Learning architecture can be effectively applied for cancer prognosis prediction with Cox-proportional hazard model incorporated. We found that Deep Learning based model demonstrated superior performances comparing to traditional machine learning models. Among the three Deep Learning based models tested, we observed that Co-net, which has the most succinct neural network structure, resulted in better prognosis performances in the measurement of concordance index and p-value of log-rank test. We showed that integrating autoencoder with Cox regression network does not significantly improve the prognosis performances."
26										"Using a combination of phenotypic, genotypic, and epigenetic parameters in glioblastoma diagnostics will bring us closer to precision medicine where therapies will be tailored to suit the genetic profile and epigenetic signature of the tumor"	Cancer progression and prognosis prediction	AUC (training + validation cohort)	external cohort validation	"Overall our study demonstrated that the Deep Learning architecture can be effectively applied for cancer prognosis prediction with Cox-proportional hazard model incorporated. We found that Deep Learning based model demonstrated superior performances comparing to traditional machine learning models. Among the three Deep Learning based models tested, we observed that Co-net, which has the most succinct neural network structure, resulted in better prognosis performances in the measurement of concordance index and p-value of log-rank test. We showed that integrating autoencoder with Cox regression network does not significantly improve the prognosis performances."
27	277 Jovčević, I.	Frontiers in Oncology	10			2020	Slovenia	https://doi.org/10.3389/fonc.2020.00281	article	Review (not applicable)	Review	Review	Review	"Using a combination of phenotypic, genotypic, and epigenetic parameters in glioblastoma diagnostics will bring us closer to precision medicine where therapies will be tailored to suit the genetic profile and epigenetic signature of the tumor"

1	Kardimalla, S and Xu, J and Link, A and Matsuyama, T and Yamamura, K and Parker, J and Uetake, H and Hernandez-Ilan, E and Lozano, J and Borazani, E and Tal, S and Evans, D and Meltzer, S J and Baba, H and Brand, R and Von Hoff, D and Balaguer, F and Lu, W and Goel, A	Cancer Research	80	16	2020	USA	1199-1206	meeting abstract	"Using this approach, we sequenced 100 plasma specimens from all GI cancers, as well as age-matched healthy controls. Eight datasets containing a total of 350 CCA, 133 adenocarcinoma and 90 HCC samples."	Case-control study	AUC (training + validation cohort)	external cohort validation	"Utilizing a novel biomarker discovery approach, we provide first evidence for cell-free DNA methylation biomarkers that offer a robust diagnostic accuracy for the identification of specific cancer types, and demonstrate their potential clinical application as a Pan-cancer panel for the early detection of all gastrointestinal cancers."
2	279 Kaur, H and Bhalra, S and Gang, D and Mehta, N and Rajghava, G P S	Journal of Oncology	73	516-517	2020	India	10.1186/s12957-020-02027-8	meeting abstract	"The dataset contains gene expression data from participants with glucose tolerance ranging from normal to newly diagnosed T2DM, in which 34 and 63 individuals were healthy and diabetic, respectively."	Case-control study	AUC, accuracy (training + validation set)	training + test set	"Prediction models developed based on three genes categorized CCA (Cholangiocarcinoma) with high precision. Thus, they can be further explored for their diagnostic and therapeutic potential for CCA."
3	280 Koshnjagt, M and Kavousi, K and Banaei-Moghaddam, A M and Moosavi-Movahedi, A A	BMC Med Genomics	13	1	2020	Iran	10.1186/s12916-020-02044-4	article	"We downloaded the FPKM-UQ (upper quartile) dataset from TCGA data portal for expression analysis"	Case-control study	AUC, accuracy, F1 score, precision, recall (10-fold CV)	cross-validation	"Using only gene expression data, it is possible to discriminate T2DM individuals from healthy controls with approximately 90 percent accuracy. Clustering of diabetic patients according to their gene expression patterns and subsequent in-depth analysis of each cluster unraveled specific abnormalities leading to insulin resistance in each cluster."
4	281 Kong, J and Lee, L and Kim, D and Han, S and Kang, H and Shin, K and Kim, S	Nat Commun	11	1	2020	Korea	10.1038/s41467-020-19313-8	article	"We separated 256 patients as the training set, 85 patients as the validation set, and 171 patients as the test set"	Case-control study	AUC (10-fold CV for training)	training + test set	"In this study, we tested if the incorporation of network analysis into an ML framework could accurately identify robust drug response biomarkers using organoid models. Indeed, we found that our method accurately predicted cancer patients' drug responses, whereas conventional ML approaches showed less optimal predictive performances. Importantly, our network-based ML model provided interpretable results for drug response prediction, which were further tested in external experimental datasets."
5	282 Koras, K and Kuravara, D and Kreis, J and Mazur, J and Staub, E and Szczerka, E	Sci Rep	10	1	2020	Poland	10.1038/s41598-020-20027-8	article	"The total set of samples consisted of 883 cancer cell lines originated from 13 tissue sites"	Drug response prediction	Correlation, RMSE (3-fold CV on training data + test set evaluation)	training + test set	"For many compounds, even a very small subset of drug-related features is highly predictive of drug sensitivity. Small feature sets selected using prior knowledge are more predictive of drug response than larger feature sets. Feature selection methods using prior knowledge to select drugs targeting specific genes and pathways, while models with wider feature sets perform better for drugs affecting general cellular mechanisms. Appropriate feature selection strategies facilitate the development of interpretable models that are indicative for therapy design."
6	283 Koureas, M and Kirgou, P and Amoutzas, G and Hadjichristodoulou, C and Gourgoulialis, K and Tsakalaki, A	Metabolites	10	8	2020	Greece	10.3390/met10081008	article	"The population sample consisted of 51 patients with confirmed LC, 38 patients with pathological computed tomography (CT) findings not diagnosed with LC, and 53 healthy controls"	Case-control study	AUC (10-fold CV + validation)	cross-validation + test set	"The random forest machine learning algorithm achieved a correct classification of patients of 88.5% (area under the curve—AUC 0.94). However, none of the methods used achieved adequate discrimination between LC patients and patients with abnormal computed tomography (CT) findings. Biomarker sets, consisting mainly of the exogenous monoaromatic compounds and 1- and 2- propanol, adequately discriminated LC patients from healthy controls."
7	284 Lai, Y H and Chen, W N and Hsu, T C and Lin, C and Tao, Y and Wu, S	Sci Rep	10	1	2020	Taiwan	10.1038/s41598-020-11596-0	article	"We accessed germline genome-wide data of 2799 early-stage breast cancer patients from the Cancer Toxicity study (NCI0203488)"	Case-control study	AUC, accuracy (10-fold CV + validation set)	training + test set	"In this study, we applied the concept of bimodal learning to construct an integrative DNN where two heterogeneous modalities (gene expression and clinical data) were integrated for predicting AUC patient overall survival. By using two modalities, the integrative DNN approach is capable of providing the missing information left by the other observed modality. Compared with our microarray DNN, we observed an increase in AUC and accuracy from the integrative DNN."
8	285 Debraze, J F and Andre, F and Vaz Luis, I	JNCI Cancer Spectrums	4	5	2020	USA	10.1093/jnci/ckaa001	article	"In this large multicentric, prospective, clinicogenomic longitudinal dataset of breast cancer survivors, we deployed machine learning techniques to investigate if high-dimensional genomic data could be used to build and validate a predictive model for the different known dimensions of fatigue. Although the ability of our models to identify clinic and genomic contributors of fatigue differed by fatigue domain, a group of SNPs and clinical variables was suggested to be associated with the cognitive domain."	Treatment response prediction	AUC (training + test set)	training + test set	"This study highlighted the potential of machine learning-aided deep methylation sequencing as a sensitive ctDNA profiling approach for early cancer detection. Further investigation in large-scale clinical studies is ongoing."
9	Li, B and Wang, C and Xu, J and Fang, S and Qiu, F and Su, J and Chu, H and Han-Zhang, H and Mao, X and Liu, H and Liu, X and Zhang, W and Zhao, H and Zhang, Z	Clinical Cancer Research	26	11	2020	China	10.1158/1078-0432.CCR.20-0167	meeting abstract	"We retained 421 MDD patients for the subsequent analysis. We performed unsupervised clustering of total 1000 HCC (hepatocellular carcinoma) samples including discovery and validation group from available public datasets"	Case-control study	AUC (training + test set)	training + test set	"In conclusion, we proposed a boosting ensemble predictive framework with the wrapper-based feature selection algorithm for predicting antidepressant treatment response and remission using an ensemble machine learning framework. The present results suggest that our boosting ensemble predictive framework with the wrapper-based feature selection algorithm may leverage a feasible way to create predictive algorithms for forecasting antidepressant treatment response and remission with clinically meaningful accuracy."
10	287 Lin, E and Kuo, P H and Liu, Y L and Yu, Y H and Yang, A C and Tsai, S J	Pharmacol Ther	13	10	2020	USA	10.1016/j.phther.2020.107605	article	"Our work demonstrated 3 immune clusters with different features. More importantly, multi-omics signatures, such as MMP9 was identified based on three clusters to help recognize patients with different prognosis and responses to immunotherapy in HCC."	Case-control study	AUC (repeated 10-fold CV)	cross-validation	"Our work demonstrated 3 immune clusters with different features. More importantly, multi-omics signatures, such as MMP9 was identified based on three clusters to help recognize patients with different prognosis and responses to immunotherapy in HCC."
11	288 Liu, F and Qin, L and Liao, Z and Song, J and Yuan, C and Liu, Y and Wang, Y and Xu, H and Zhang, Q and Pei, Y and Zhang, H and Pan, Y and Chen, X and Zhang, Z and Zhang, W and Zhang, B	Hematology and Oncology	9	1	2020	China	10.1016/j.hbon.2020.100163	article	"The discovery stage involved 160 pairs of ccRCC (clear cell renal cell carcinoma) and matched normal tissues for investigation of DNAm and biomarkers as well as 318 cases of ccRCC including clinical signatures"	Prognostic subtype stratification	correlation (discovery + validation cohorts)	external cohort validation	"Our work demonstrated 3 immune clusters with different features. More importantly, multi-omics signatures, such as MMP9 was identified based on three clusters to help recognize patients with different prognosis and responses to immunotherapy in HCC."
12	289 Liu, P and Tian, W	PeerJ	8		2020	China	10.7717/peerj.9624	article	"The present study provides a comprehensive analysis of ccRCC using multi-omics data. These findings indicated that multi-omics analysis could identify some novel epigenetic factors, which were the most important causes of advanced cancer and poor clinical prognosis."	Case-control study	AUC (10-fold CV + validation cohort)	cross-validation + external cohort validation	"The present study provides a comprehensive analysis of ccRCC using multi-omics data. These findings indicated that multi-omics analysis could identify some novel epigenetic factors, which were the most important causes of advanced cancer and poor clinical prognosis."
13	290 Liu, X Y and Wu, S B and Zeng, W Q and Yuan, J J and Xu, H B	Sci Rep	10	1	2020	China	10.1038/s41598-020-18180-9	article	"This data describes 20,501 genes in 806 different breast cancer samples. We retained only samples with complete information. After that, 85 TNBC and 460 non-TNBC were further divided into two groups: training (n = 327; 51 TNBC, 276 non-TNBC) and testing (n = 218; 34 TNBC, 184 non-TNBC) sets"	Case-control study	AUC (10-fold CV + validation)	cross-validation + test set	"The present study provides a comprehensive analysis of ccRCC using multi-omics data. These findings indicated that multi-omics analysis could identify some novel epigenetic factors, which were the most important causes of advanced cancer and poor clinical prognosis."
14	291 Liu, Y and Liu, F and Hu, X and He, J and Jiang, Y	Journal of Oncology	14		2020	China	10.1186/s12957-020-02060-0	article	"The most significant contribution of this article is the integration of the genetic mutation and expression profiles to determine prognostic genes for LIAD patients. If genetic expression and mutation profiles are available, the pipelines of determining DEGs and DMGs in this article can be applied to other types of cancers."	Prognostic study	AUC (10-fold CV + validation)	cross-validation + test set	"The most significant contribution of this article is the integration of the genetic mutation and expression profiles to determine prognostic genes for LIAD patients. If genetic expression and mutation profiles are available, the pipelines of determining DEGs and DMGs in this article can be applied to other types of cancers."
15	292 Lu, Y and Wu, S and Gu, C and Wu, M and Wang, S and Yue, Y and Liu, M and Sun, Z	OncoTarget and Therapy	13		2020	China	10.21601/ott.57955	article	"Our research results indicate that the 8 gene prognostic signature is a reliable tool for predicting the OS (overall survival) of CMMI (colon adenocarcinoma) patients."	Prognostic study	AUC (training + test set + external validation set)	external cohort validation	"In this study, the AUC of 8 gene signature screened by multi-omics in the training set and validation set for five years is more than 0.64, which is more effective in predicting the prognosis of patients."
16	293 Luca, B A and Moulton, V and Ellis, C and Edwards, D R and Campbell, C and Cooper, R A and Clark, J and Bressan, D S and Cooper, C S	Br J Cancer	122	10	2020	Kingdom	10.1038/s41416-020-02061-8	article	"There were 1785 samples from primary malignant tissue, and 173 from normal tissue. A total of 901 TCGA NSCLC samples were available using the Illumina Infinium HumanMethylationEP480 platform, including 827 tumor tissues and 74 non-tumor tissues"	Prognostic subtype stratification	Correlation, log-rank p-value (hold-out validation)	training + test set	"We have confirmed a key prediction of the DESNT cancer model by demonstrating that the presence of a small proportion of the DESNT cancer signature confers poorer outcome. The proportion of DESNT signature can be considered a continuous variable, such as DESNT cancer content increases, the outcome became worse. This observation led to the development of nomograms for estimating PSA failure at 3, 5 and 7 years following prostatectomy."
17	294 Luo, R and Song, J and Xiao, X and Xiao, Z and Zhao, Z and Zhang, W and Miao, S and Tang, Y and Ran, L	Aging (Albany NY)	12	14	2020	China	10.1080/10791808.2020.1817000	article	"This included 872 samples that had sequenced genomes, clinical data, and 50K normalized transcripts from RNA sequencing"	Tumor recurrence and immunotherapy response prediction	AUC (training + external validation)	external cohort validation	"In this study, we initially identified 4 CpG biomarkers associated with recurrence of NSCLC. Based on TCGA NSCLC cohort comprised of lung adenocarcinoma (LUAD) and lung squamous cell carcinoma (LUSC), a promising DNAm-based risk score model predictive of disease was constructed and then validated in the other 3 datasets."
18	295 Makarios, M and Iwaki, H and Blauwendraat, C and Leonard, H and Hashemi, S and Kim, J and Van Keuren-Jensen, K and Craig, D and Appelmann, E and Smolevsky, I and Bookman, M and Singleton, A and Fagnini, F and Naik, M	Movement Disorders	35		2020	USA	10.1080/15402002.2020.1818000	meeting abstract	"Machine learning is a promising approach in the search of more accurate and generalizable models for prediction of CAT (Cancer-Associated Venous Thromboembolism). In the application described here, the use of random survival forests performed well without information about future chemotherapy administration. Additional work is needed to identify the optimal algorithm and covariates, including better delineation of which cancer genomic information should be retained."	Case-control study	AUC (30K test samples after training on 70% of samples)	training + test set	"Machine learning is a promising approach in the search of more accurate and generalizable models for prediction of CAT (Cancer-Associated Venous Thromboembolism). In the application described here, the use of random survival forests performed well without information about future chemotherapy administration. Additional work is needed to identify the optimal algorithm and covariates, including better delineation of which cancer genomic information should be retained."
19	296 Mantha, S and Dunbar, A and Bolton, K L and Devlin, S and Gorenshvyn, D and Donoghue, M and Arora, M E and Soff, G A	Blood	136		2020	USA	10.1182/blood.2020.118829	meeting abstract	"12,040 patients were included in the final analysis. There were 855 CAT events during the observation period"	Prognostic study	C-index (cross-validation)	cross-validation	"Multivariate models based on support vector machines and the LASSO variable selection method yielded two discriminant lipid panels for cRCC detection and early diagnosis. A lipid panel allowed discriminating cRCC patients from controls with 95.7% accuracy in training set under cross-validation and 77.1% accuracy in an independent test set."
20	297 Manzi, M and Palazzo, M and Kott, M E and Beausery, P and Kankiewicz, P and Giménez, M and Monge, M E	J Proteome Res	20	1	2021	Argentina	10.1021/acs.jproteome.1c00063	article	"In a derivation cohort of 637 patients referred for coronary angiography, predictors of 270% coronary stenosis were identified from 6 clinical variables and 109 biomarkers. The final model was first internally validated on a separate case (n=273) and then externally validated on a cohort of 241 patients"	Case-control study	accuracy (training/test set)	training + test set	"We have derived and externally validated a clinical/proteomic panel that can predict the presence of obstructive CAD (Coronary Artery Disease) with high accuracy. The score performs similarly well in the evaluation of acute chest pain in the ED (including in-patients referred in non-acute chest pain) and in outpatients presenting for evaluation of stable angina including those with renal injury."
21	298 McCarthy, C P and Neumann, J T and Michelunaso, S A and Ibrahim, N E and Gagnan, H K and Siemsen, W A and Schläpfer, S and Miller, T and Maguire, C A and Barnes, C and Rhyne, K F and Westermann, D and Januzzi, Jr, J L	J Am Heart Assoc	9		2020	USA	10.1161/JAHA.120.02211	article	"We have derived and externally validated a clinical/proteomic panel that can predict the presence of obstructive CAD with high accuracy. The score performs similarly well in the evaluation of acute chest pain in the ED (including in-patients referred in non-acute chest pain) and in outpatients presenting for evaluation of stable angina including those with renal injury."	Case-control study	AUC (train + test + external validation)	external cohort validation	"We have derived and externally validated a clinical/proteomic panel that can predict the presence of obstructive CAD with high accuracy. The score performs similarly well in the evaluation of acute chest pain in the ED (including in-patients referred in non-acute chest pain) and in outpatients presenting for evaluation of stable angina including those with renal injury."

291	Miao, R and Chen, H H and Dang, Q X and Li, Y and Yang, Z Y and He, M F and Xiao, Z F and Liang, Y	Pharmacol Res	159	104932-104932	2020	China	https://doi.org/10.1016/j.phrs.2020.104932	article	"The GDSC dataset contains 140 drug sensitivity experiments results in 624 cell lines." "Targeted DNA sequencing for more than 500 cancer-associated genes and exome-capture RNA sequencing was carried out in more than 25,000 fresh frozen or paraffin embedded tumor samples, including both primary and metastatic tumors." "We divided 741 ADNI participants with blood microarray data into three groups based on their most recent CDR assessment: cognitive normal (CDR = 0), mild cognitive impairment (CDR = 0-3), and probable Alzheimer's disease (CDR < 1.0)." "Here, we investigate whether 22 VOCs from the breath of 246 patients can distinguish those with no liver disease (n = 54), cirrhosis (n = 30), HCC (n = 112), pulmonary hypertension (n = 49), or colorectal cancer liver metastases (n = 51)."	Drug response prediction	AUC, sensitivity, specificity (5-fold CV)	cross-validation	"The proposed model of this paper used statistical methods and Machine Learning methods combined with genomic data to accurately predict the performance of oncology drugs on cancer cell lines." "The proposed model of this paper used statistical methods and Machine Learning methods combined with genomic data to accurately predict the performance of oncology drugs on cancer cell lines."	
300	T Michuda, J and Ielbowitz, B and Aman-Farash, S and Bevis, C and Breschi, A and Kapilovsky, J and Igitau, C and Bell, J S and Bouchamp, K A and White, K and Stumpe, M and Beaulieu, M and Taster, R	Cancer	80	16	2020	USA	https://doi.org/10.1186/s12864-020-05423-4	meeting abstract	"Multimodal prediction of diagnosis for cancers of unknown primary"	Differential diagnosis prediction	accuracy (training/test set)	training + test set	"The incorporation of multiple modes of omics data can improve the interpretability and robustness of machine learning models to predict cancer diagnosis" "Our analyses indicate that machine learning may be able to predict cognitive decline in individuals using RNA levels from a blood microarray by taking into account small differences in expression that are individually nonsignificant. A support vector machine was able to increase predictive accuracy of AD from a 55% baseline to almost 90%."	
301	Miller, J B and Kauwe, J S K	Genes (Base)	11	6	2020	USA	https://doi.org/10.1007/s11065-020-09572-0	article	"Predicting Clinical Dementia Rating Using Blood RNA Levels"	Differential diagnosis prediction	AUC (10-fold CV)	cross-validation	"Our analyses indicate that machine learning may be able to predict cognitive decline in individuals using RNA levels from a blood microarray by taking into account small differences in expression that are individually nonsignificant. A support vector machine was able to increase predictive accuracy of AD from a 55% baseline to almost 90%."	
302	Miller-Atkins, G and Acevedo-Moreno, L A and Grove, D and Dweik, R A and Tomelli, A R and Brown, J M and Allende, D S and Aucejo, F and Rotroff, D M	Hepatology Communications	4	7	1041-1055	2020	USA	https://doi.org/10.1007/s40611-020-01490-9	article	"Breath Metabolomics Provides an Accurate and Noninvasive Approach for Screening Cirrhosis, Primary, and Secondary Liver Tumors"	Differential diagnosis prediction	balanced accuracy (cross-validation)	cross-validation	"The use of machine learning and breath VOCs (volatile organic compounds) shows promise as an approach to develop improved, noninvasive screening tools for chronic liver disease (primary and secondary liver tumors)."
303	Mongan, D and Focking, M and Healy, C and Susal, R and Cagney, G and Cannon, M and Zammitt, S and Nelson, B and McGorry, P and Nordentoft, M and Krebs, M G and Riecher-Rossler, A and Bressan, R and Barrettes-Vidal, N and Borgwardt, S and Rahmsdorf, S and Sachs, G and Van Der Gaag, M and Rutten, B A and Fostel, C and De Haan, L and Valmianski, L and Kempton, M and McGuire, P and Cotter, D	Schizophrenia Bulletin	46	6	5238-5239	2020	Ireland	https://doi.org/10.1093/schbul/kbaa011	meeting abstract	"Development of proteomic prediction models for outcomes in the clinical high risk state and psychotic experiences in adolescence: Machine learning analysis of nested case-control studies"	Case-control study	AUC, PPV, NPV (training + test set)	training + test set	"With external validation, models incorporating proteomic data may contribute to improved prediction of clinical outcomes in individuals at risk of psychosis" "Taken together, we have presented three unique CNN architectures that take high dimension gene expression inputs and perform cancer type prediction while considering their tissue of origin. Our model achieved an equivalent 95.7% prediction accuracy compared to earlier published studies, however with a drastically simplified CNN construction and with a reduced influence of the tissue origins."
304	Mostaf, M and Chi, Y C and Huang, Y and Chen, Y	BMC Med	13	44-44	2020	USA	https://doi.org/10.1186/s12916-020-02002-0	article	"Convolutional neural network models for cancer type prediction based on gene expression"	Differential diagnosis prediction	accuracy (5x 6-fold CV, 80-20% splitting for training and validation)	cross-validation	"Pharmacogenomic biomarkers including gene variants for cancer susceptibility genes (CASC15) and important MTX pathway enzymes (ATIC) combined with baseline DAS28 score predicted MTX response in patients with early RA more reliably than demographic and baseline DAS28 alone, with replication in an independent cohort"	
305	Myasoedova, E and Athreya, A and Crowson, C and Weisshilbourn, R and Wang, L and Matteson, E T	Arthritis and Rheumatism	72	4014-4025	2020	USA	https://doi.org/10.1002/art.41738	meeting abstract	"Individualized Prediction of Response to Methotrexate Treatment in Patients with Rheumatoid Arthritis: A Pharmacogenomic-driven Machine Learning Approach"	Drug response prediction	AUC (5x 10-fold CV + external validation)	cross-validation + external cohort validation	"Pharmacogenomic biomarkers including gene variants for cancer susceptibility genes (CASC15) and important MTX pathway enzymes (ATIC) combined with baseline DAS28 score predicted MTX response in patients with early RA more reliably than demographic and baseline DAS28 alone, with replication in an independent cohort"	
306	Naz, H and Ahuja, S	Diabetes and Metabolic Disorders	19	1	391-403	2020	India	https://doi.org/10.1007/s00125-020-05005-8	article	"Deep learning approach for diabetes prediction using PIMA indian dataset"	Case-control study	accuracy ("Splits in an 80/20% ratio into the training and validation set")	training + test set	"The outcomes of this study confirms that DL provides the best results with the most promising extracted features. DL achieves the accuracy of 98.07% which can be used for further development of the automatic prognosis tool."
307	Nasha, A and Sekeres, M A and Bejar, R and Rauh, M J and Othum, M and Komroji, R S and Barnard, J and Hilton, C B and Kerr, C M and Stensson, D P and DeZern, A and Roboz, G and Garcia-Manero, G	JCO Precision Oncology	3	2019	USA		https://doi.org/10.1200/JCO.2019.37.111	article	"Genomic biomarkers to predict response to hypomethylating agents in patients with myelodysplastic syndromes using artificial intelligence"	Drug response prediction	accuracy (training/test set)	training + test set	"Genomic biomarkers can identify, with high accuracy, approximately one third of patients with MDS who will not respond to HMAs. This study highlights the importance of machine learning technologies such as the recommender system algorithm in translating genomic data into useful clinical tools."	
308	Nielsen, R L and Helemlus, M and Garcia, S L and Roager, H M and Aytan-Aktug, D and Hansen, L B S and Lind, M Y and Vogt, J K and Dalgaard, M D and Bahl, M I and Jensen, C B and Mukhopadhyay, R and Wainroe, C and Aakjov, V and Gabel, R and Riström, M and Falke, H and Sparholt, M H and Christensen, A F and Vestergaard, H and Hansen, T and Kristiansen, K and Erik, S and Petersen, T N and Lauritzen, L and Licht, T R and Pedersen, G and Gupta, R	Sci Rep	10	1	20193-20193	2020	Denmark	https://doi.org/10.1038/s41598-020-70072-7	article	"Here, we classify weight loss responders (N = 106) and non-responders (N = 97) of overweight non-diabetic middle-aged Danes to two earlier reported dietary trials over 8 weeks"	Treatment response prediction	AUC ("50 shuffle-split fivefold cross-validations was used")	cross-validation	"By identifying the propensity of study participants likely to experience weight loss, a more effective individual targeting of dietary interventions can be facilitated, eventually in concert with comprehensive population weight loss strategies. Furthermore, understanding predictive features of weight loss responses will drive improved understanding of the interplay between gut microbiota, diet and individual predisposition."
309	Nyamundanda, G and Eason, K and Guiney, J and Lord, C J and Sadaanandam, A	Cancers	12	10	1-14	2020	United Kingdom	https://doi.org/10.3390/cancers12110241	article	"In total, 2043 breast cancer samples were used in this work"	Subgroup stratification	Silhouette width, cophenetic correlation (external test datasets)	external cohort validation	"Overall, this genome-phenome machine-learning integration tool, PhenMap identifies functional and phenotype-integrated discrete or continuous subtypes with clinical translational potential."
310	Ozer, M E and Sarica, P O and Aga, Y K	Omics	24	5	241-246	2020	Turkey	https://doi.org/10.3390/omics24050241	article	"New Machine Learning Applications to Accelerate Personalized Medicine in Breast Cancer: Rise of the Support Vector Machines"	review (not applicable)	Review	review (not applicable)	"This expert review describes and examines, first, the SVM models employed to forecast breast cancer subtypes using diverse systems science data, including transcriptomics, genomics, proteomics, and radiomics, as well as biological pathway, clinical, pathological, and biochemical data. Then, we compare the performance of the present SVM and other diagnostic and therapeutic prediction models across the data types. We conclude by emphasizing that data integration is a critical bottleneck in systems science, cancer research and development, and health care innovation and that SVM and machine learning approaches offer new solutions and ways forward in biomedical, bioengineering, and clinical applications."
311	Pal, S and Weber, P and Isserlin, A and Kalka, H and Hui, S and Shah, M A and Giudice, L and Giugno, R and Nèp, A and Baumbach, J and Bader, G D	F1000res	9	1239-1239	2020	Canada	https://doi.org/10.1093/f1000/2020/09/1239	article	"Including 154 Luminal A and 194 tumours of other subtypes" "We conduct a variety of simulations and trials against the Madelon benchmark dataset from the University of California – Irvine (UCI), and Clinical data from The Cancer Genome Atlas (TCGA)."	Case-control study	AUROC, AUPR, and accuracy (an approximately 70:30 split of samples was used for cross validation)	cross-validation	"The netBx Bioconductor package provides a novel workflow for pathway-based patient classification from sparse genetic data." "biomarkers identified upon previous methods for identifying interpretable features in RFS and bring them together under a correlated binomial distribution to create an efficient hypothesis testing algorithm that identifies biomarkers' main effects and interactions. Preliminary results in simulations demonstrate computational gains while retaining competitive model selection and classification accuracies."	
312	Rashid Zaim, S and Kenost, C and Berghout, J and Chi, W and Wilson, L and Zhang, H H and Lussier, Y	BMC Bioinformatics	21	1	374-374	2020	USA	https://doi.org/10.1186/s12859-020-03713-8	article	"Clinical and genomic data, including commercially available next generation sequencing panels, were obtained for patients (91) treated at the Cleveland Clinic (CC; 63 pts), Munich Leukemia Laboratory (ML; 1509 pts), and the University of Pavia in Italy (UP; 536 pts)."	Case-control study	Precision, recall, test error (training and test set)	training + test set	"We developed and externally validated a highly accurate and interpretable model that can distinguish MDS from other myeloid malignancies using clinical and mutational data from a large international cohort. The model can provide personalized interpretations of its outcome and can aid physicians and hematopathologists in recognizing MDS with high accuracy when encountering pts with pancytopenia and with a suspected diagnosis of MDS."
313	Radokovich, N and Meggersdorfer, M and Malcovati, L and Sekeres, M A and Shreve, J and Beau Hilson, C and Roushali, Y and Walker, W and Hutter, S and Mukherjee, S and Kerr, C M and Jhu, B K and Gall, A and Pozzi, S and Gerds, A T and Hiferlach, C and Maciejewski, J P and Hiferlach, T and Nizha, A	Blood	136	33-35	2020	USA	https://doi.org/10.1182/blood-2020-130412	meeting abstract	"A personalized clinical-decision tool to improve the diagnostic accuracy of myelodysplastic syndromes"	Case-control study	AUC (training + external validation)	external cohort validation	"In this article, we compare the usefulness and limitations of traditional statistical methods and ML, when applied to the medical field. Traditional statistical methods seem to be more useful when the number of cases largely exceeds the number of variables under study and a priori knowledge on the topic under study is substantial such as in public health. ML could be more suited in highly innovative fields with a huge bulk of data, such as omics, radiogenomics, drug development, and personalized treatment. Integration of the two approaches should be preferred over a unidirectional choice of either approach."	
314	Rajula, H S R and Verlotto, G and Manchia, M and Antonucci, N and Fanos, V	Medicina	56	9	945-945	2020	Italy	https://doi.org/10.3390/med5609045	article	"Comparison of Conventional Statistical Methods with Machine Learning in Medicine: Diagnosis, Drug Development, and Treatment"	review (not applicable)	Review	review (not applicable)	"We have developed a DNA methylation score for exposure to maternal smoking during pregnancy, outperforming the three previously developed scores. One possible application of the current score could be for model adjustment purposes or to assess its association with distal health outcomes where part of the effect can be attributed to maternal smoking. Further, it may provide a biomarker for fetal exposure to maternal smoking."
315	Rauscher, S and Molton, P E and Heiskala, A and Karhunen, V and Burdige, G and Craig, J M and Godfrey, K M and Lillyrop, K and Mori, T A and Bell, L J and Oddy, W H and Proulx, C and Jewell, M R and Sebert, S and Huang, R C	Environ Health Perspect	128	9	97003-97003	2020	Australia	https://doi.org/10.1289/ehp.201097003	article	"The Raine study was developed and tested in the Raine Study with data from 995 white 17-y old participants using 10-fold cross-validation"	Case-control study	Sensitivity, specificity (10-fold CV)	cross-validation	"In this research, we compared three machine learning methods that have been proved to construct powerful predictive models (genetic algorithms, LASSO, and stepwise) and propose the inclusion of markers from misclassified samples to improve overall prediction accuracy. Our results show that the addition of markers from an initial model and the markers of the model fitted to misclassified samples improves the area under the receiving operative curve by around 5%, reaching ~0.84, which is highly competitive using only genetic information"
316	Ristori, M V and Mortera, S I and Marzano, V and Guerrero, S and Vermocchi, P and Iaino, G and Gardini, S and Torricelli, G and Valeri, G and Vicari, S and Gabarrini, A and Putignano, L	Int J Mol Sci	21	17		2020	Italy	https://doi.org/10.3390/ijms21170241	article	"The study has four groups accounting for 4,230 individuals"	Case-control study	AUC (20 rounds of internal cross-validation (CV) to 80% of the dataset for training and 20% for testing)	cross-validation	"In this research, we compared three machine learning methods that have been proved to construct powerful predictive models (genetic algorithms, LASSO, and stepwise) and propose the inclusion of markers from misclassified samples to improve overall prediction accuracy. Our results show that the addition of markers from an initial model and the markers of the model fitted to misclassified samples improves the area under the receiving operative curve by around 5%, reaching ~0.84, which is highly competitive using only genetic information"
317	Romero-Rosales, B L and Tames-Pena, J G and Nicolini, H and Moreno-Treviño, M G and Treviño, V	PLoS One	15	4	e0232103-0232103	2020	Mexico	https://doi.org/10.1371/journal.pone.0232103	article	"The raw data from 13 synovial datasets with 284 samples used in 4 blood datasets with 1,885 samples were downloaded and processed"	Case-control study	AUC (training + test set)	training + test set	"The raw data from 13 synovial datasets with 284 samples used in 4 blood datasets with 1,885 samples were downloaded and processed"
318	Rychlik, D and Neely, J and Sirota, M	Arthritis and Rheumatism	72	1503-1504	2020	USA	https://doi.org/10.1002/art.41738	meeting abstract	"Uncovering Novel Biomarkers for Rheumatoid Arthritis from Feature Selection Approaches on Synovium and Blood Gene Expression Data"	Case-control study	AUC (training + test set)	training + test set	"This novel list of biomarkers, identified through a robust feature selection procedure on public data and validated using multiple independent data sets, coupled with the RAIScore may be useful in the early diagnosis and disease and treatment monitoring of RA."	

