



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

Empirical evidence on the impact of study characteristics and the performance of prediction models: a meta-epidemiological study

Journal:	<i>BMJ Open</i>
Manuscript ID	bmjopen-2018-026160
Article Type:	Research
Date Submitted by the Author:	22-Aug-2018
Complete List of Authors:	Damen, Johanna; Cochrane Netherlands, University Medical Center Utrecht, Utrecht University; Julius Center for Health Sciences and Primary Care, University Medical Center Utrecht, Utrecht University Debray, Thomas; Julius Center for Health Sciences and Primary Care, University Medical Center Utrecht, Utrecht University; Cochrane Netherlands, University Medical Center Utrecht, Utrecht University Pajouheshnia, Romin; Julius Center for Health Sciences and Primary Care, University Medical Center Utrecht, Utrecht University Reitsma, Johannes; Julius Center for Health Sciences and Primary Care, University Medical Center Utrecht, Utrecht University; Cochrane Netherlands, University Medical Center Utrecht, Utrecht University Scholten, Rob; Cochrane Netherlands, University Medical Center Utrecht, Utrecht University; Julius Center for Health Sciences and Primary Care, University Medical Center Utrecht, Utrecht University Moons, Karel; Julius Center for Health Sciences and Primary Care, University Medical Center Utrecht, Utrecht University; Cochrane Netherlands, University Medical Center Utrecht, Utrecht University Hooft, Lotty; Cochrane Netherlands, University Medical Center Utrecht, Utrecht University; Julius Center for Health Sciences and Primary Care, University Medical Center Utrecht, Utrecht University
Keywords:	Meta-epidemiology, Prognosis, Prognostic models, Bias, Prediction

SCHOLARONE™
Manuscripts

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

1 Empirical evidence on the impact of study characteristics and the performance of prediction models: a
2 meta-epidemiological study

3
4
5
6
7
8 Johanna A A G Damen, Thomas P A Debray, Romin Pajouheshnia, Johannes B Reitsma, Rob J P M
9 Scholten, Karel G M Moons, Lotty Hooft
10
11
12
13
14
15 Johanna A A G Damen
16 Assistant professor
17
18 Cochrane Netherlands, University Medical Center Utrecht, Utrecht University, Utrecht, The Netherlands;
19 Julius Center for Health Sciences and Primary Care, University Medical Center Utrecht, Utrecht
20 University, 3508 GA Utrecht, The Netherlands
21
22
23
24
25 Thomas P A Debray
26 Assistant professor
27
28 Cochrane Netherlands, University Medical Center Utrecht, Utrecht University, Utrecht, The Netherlands;
29 Julius Center for Health Sciences and Primary Care, University Medical Center Utrecht, Utrecht
30 University, 3508 GA Utrecht, The Netherlands
31
32
33
34
35 Romin Pajouheshnia
36 PhD fellow
37
38 Julius Center for Health Sciences and Primary Care, University Medical Center Utrecht, Utrecht
39 University, 3508 GA Utrecht, The Netherlands
40
41
42
43 Johannes B Reitsma
44 Associate professor
45
46 Cochrane Netherlands, University Medical Center Utrecht, Utrecht University, Utrecht, The Netherlands;
47 Julius Center for Health Sciences and Primary Care, University Medical Center Utrecht, Utrecht
48 University, 3508 GA Utrecht, The Netherlands
49
50
51
52
53 Rob J P M Scholten
54 Professor
55
56
57
58
59
60

33 Cochrane Netherlands, University Medical Center Utrecht, Utrecht University, Utrecht, The Netherlands;
34 Julius Center for Health Sciences and Primary Care, University Medical Center Utrecht, Utrecht
35 University, 3508 GA Utrecht, The Netherlands

37 Karel G M Moons

38 Professor

39 Cochrane Netherlands, University Medical Center Utrecht, Utrecht University, Utrecht, The Netherlands;
40 Julius Center for Health Sciences and Primary Care, University Medical Center Utrecht, Utrecht
41 University, 3508 GA Utrecht, The Netherlands

43 Lotty Hooft

44 Associate professor

45 Cochrane Netherlands, University Medical Center Utrecht, Utrecht University, Utrecht, The Netherlands;
46 Julius Center for Health Sciences and Primary Care, University Medical Center Utrecht, Utrecht
47 University, 3508 GA Utrecht, The Netherlands

49 Correspondence to:

50 Johanna A A G Damen

51 Julius Center for Health Sciences and Primary Care, University Medical Center Utrecht

52 P.O. Box 85500

53 Str. 6.131

54 3508 GA Utrecht

55 The Netherlands

56 j.a.a.damen@umcutrecht.nl

57 +31 88 75 693 77

58

59 Word count: 4167

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

Objectives: To empirically assess the relation between study characteristics and prognostic model performance in external validation studies of multivariable prognostic models.

Design: Meta-epidemiological study.

Data sources and study selection: We searched electronic databases for systematic reviews of prognostic models published in the period 2010-2016. Reviews from non-overlapping clinical fields were selected if they reported common performance measures (either the concordance (c)-statistic or the ratio of observed over expected number of events (OE ratio)) from ten or more validations of the same prognostic model.

Data extraction and analyses: Study design features, population characteristics, methods of predictor and outcome assessment, and the aforementioned performance measures were extracted from the included external validation studies. Random effects meta-regression was used to quantify the association between the study characteristics and model performance.

Results: We included 10 systematic reviews, describing a total of 224 external validations, of which 221 reported c-statistics and 124 OE ratios. Associations between study characteristics and model performance were heterogeneous across systematic reviews. C-statistics were most associated with variation in population characteristics, outcome definitions and measurement, and predictor substitution. For example, validations with eligibility criteria comparable to the development study were associated with higher c-statistics compared to narrower criteria (difference in logit c-statistic 0.21 [95% CI 0.07, 0.35], similar to an increase from 0.70 to 0.74). Using a case-control design was associated with higher OE ratios, compared to using data from a cohort (difference in log OE ratio 0.97 [95% CI 0.38, 1.55], similar to an increase in OE ratio from 1.00 to 2.63).

Conclusions: Variation in performance of prognostic models across studies is mainly associated with variation in case-mix, study designs, outcome definitions and measurement methods, and predictor substitution. Researchers developing and validating prognostic models should realise the potential influence of these study characteristics on the predictive performance of prognostic models.

Strengths and limitations of this study

- To the best of our knowledge, this is the first meta-epidemiological study focusing on the association of study characteristics with estimates of prognostic model performance.
- We included all ten systematic reviews describing at least ten external validations of the same prognostic model, resulting in 224 external validations.
- We extracted relevant features of design and conduct according to existing checklists on quality assessment (CHARMS) and reporting of prediction model studies (TRIPOD).
- It was not feasible to fit multivariable meta-regression models due to the limited number of available, well-reported, validation studies within the individual reviews, rendering the effective sample size too small for multivariable meta-regression analyses.

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

Introduction

Prediction models, including diagnostic and prognostic models, estimate the probability that an individual has or will develop a certain outcome (e.g. disease or complication). Hereto, they combine multiple predictors into an estimate of an individual's risk.¹ Before using a prediction model in clinical practice it is recommended to validate the performance of the model in a population other than the population in which the model was developed (so called external validation studies).² Such studies assess whether model predictions remain sufficiently accurate across different settings and populations. Obviously, it is important that the methodological quality of external validation studies is good, as otherwise estimates of the prediction model's performance may be biased and thereby lead to misleading conclusions on its generalizability to practice.

Systematic reviews have found that the performance of existing prediction models often varies substantially across external validation studies of those models.³⁻⁵ These differences may not only appear due to random variation (when validation studies are small), but may also arise when model predictions are invalid because the model is applied in very different populations (eg, the association between predictors in the model and the outcome are different) or when design-related characteristics of the validation study (eg, measurement methods or variable definitions) are not well aligned with the original development study.^{2,6}

To provide empirical evidence of the association of study characteristics with prediction model performance, a meta-epidemiological approach can be used. Studies using this approach have shown the influence of study characteristics on the effectiveness of interventions studied in randomized trials and on the accuracy of diagnostic tests.⁷⁻¹² For diagnostic prediction models evidence suggests estimates of performance may be biased in studies with certain study characteristics. One study found a higher diagnostic odds ratio in case-control studies, studies with differential outcome verification (ie, using different outcome assessments across study individuals), and with low sample size.¹³ To date, no meta-epidemiological study has been performed investigating the possible impact of study characteristics on measures of the predictive performance of a prognostic model upon external validation, which is commonly quantified in terms of discrimination and calibration.¹⁴ The aim of this study was to investigate sources of heterogeneity in the predictive performance of prognostic models. A meta-

127 epidemiological approach was used to synthesize evidence from a range of clinical fields. This study can
128 serve as empirical evidence for design and analysis related bias in prognostic model studies.

For peer review only

Methods

Search and selection of systematic reviews

We used an existing database (last updated on March 27, 2017) consisting of studies evaluating multiple existing prediction models, including narrative or systematic reviews of prediction models, or head-to-head comparisons of multiple prediction models validated on a specific dataset (See Supplement 1 for details of the search strategy). To construct this database, references identified by the search were screened for eligibility by one reviewer (GSC) on title, abstract and, if necessary, on full text. Subsequently, the full text of all articles in the database were screened for eligibility to the current project by another reviewer (JAAGD). We selected systematic reviews of prognostic models (ie, diagnostic models were excluded) that included at least ten studies that externally validated the same prognostic model, and that presented the performance of these models in terms of discrimination (concordance (c)-statistic or area under the receiver operating characteristic (AUC) curve), or calibration (observed expected (OE) ratio). Discrimination is the ability of the model to distinguish between people who will and who will not develop the outcome of interest, while calibration reflects the overall agreement between the total number of observed and predicted ('expected') events.¹⁴ We excluded systematic reviews that selected studies based on specific study characteristics (eg, we excluded systematic reviews that did not include primary studies with a sample size below 100, if we were not able to identify the primary studies that had been excluded for this reason). Furthermore, we excluded reviews of prognostic models in which the weights of predictors in the original model were based on expert opinion rather than on coefficients estimated from a formal statistical approach. If more than one systematic review on the same prognostic model was identified, we included the one with the broadest inclusion criteria (eg, reviews focussing on specific patient populations were not preferred if a review with a broader population was available) or the most recent review (in this order of preference). When multiple prognostic models for the same condition were described in one systematic review which all fulfilled the selection criteria, we included the model with the highest number of external validations.

Selection of the primary external validation studies from the included systematic reviews

From the included systematic reviews we collected the primary studies in which the prognostic models were developed and externally validated. For primary external validation studies for which no measure

of discrimination (c-statistic) or calibration (OE ratio) was reported in the systematic review, we checked the full text of the primary study, and if performance was not reported, these studies were excluded. If primary external validation studies described multiple external validations of the same model and if there was no overlap in included participants between these external validations (eg, a model was validated in two different cohorts, or a model was validated in men and women separately), data were extracted for every external validation separately. If a model was validated multiple times on the same population (described in either one or multiple publications), we selected the external validation that was included in the systematic review. If the systematic review included all those external validations, we selected the one in which the study population and predicted outcome most closely resembled the population and outcome of the original model.

Data extraction and preparation

We extracted relevant features of design and conduct according to existing checklists on quality assessment (CHARMS) and reporting of prediction model studies (TRIPOD).¹⁵⁻¹⁷ Information about study characteristics of studies in which the models were developed were extracted from the corresponding development papers. Information about study characteristics of primary external validation studies were first extracted from systematic reviews. This information was subsequently checked using the external validation studies and, if necessary, additional information was extracted by one reviewer (JAAGD or RP). Items we extracted included study type (eg, external validation only, development of a new model and external validation of a model), study design (eg, existing cohort, existing RCT), dependency of investigators (validation by independent investigators or investigators also involved in the development study), eligibility criteria for participant inclusion, setting, location (continent), study dates, number of centres, follow-up time and prediction horizon, age and gender distribution, deletion or substitution of predictors, outcome definition and measurement method, sample size and number of events, handling of missing data, and model performance (see Supplement 2 for details). The data extraction form was piloted on multiple articles by all reviewers (JAAGD, TPAD, LH, KGMM, RP, JBR, RJPMS).

For analysis purposes, some study characteristics had to be categorized or transformed (Supplement 2). For example, eligibility criteria of the validation study as compared to the development study had to be judged and categorized as comparable, narrower (if subgroups included in the development study were excluded from the validation study), broader (if subgroups excluded from the development study were included in the validation study), mixture (a combination of the two), or unclear. For setting, location,

predictors and outcome a similar categorization was used. If data on study characteristics were not reported in the primary external validation studies, these were either categorized as ‘unclear’ (in case of categorical study variables), or the study was excluded from the analyses of that (missing) study characteristic (in case of continuous study variables, such as sample size). In order to improve comparability between reviews, we standardized continuous study variables separately for each systematic review, i.e. for every variable we subtracted the mean and divided by the standard deviation of all external validations identified from the same systematic review.

Statistical analyses

We used a two-staged approach to study the possible association between study characteristics and predictive performance.

In the first stage, we fitted a univariable meta-regression model for every study characteristic within each systematic review with the logit c-statistic or log OE ratio as outcome variable.¹⁸ The regression coefficients estimated from this meta-regression model indicate the difference in logit c-statistic or log OE ratio between a certain category of a study characteristic and a chosen reference category (ie, the category that was present in most systematic reviews) of that characteristic.

In the second stage, these regression coefficients were pooled by the use of a random effects model. This reflected the average influence of the study characteristic on model performance across all systematic reviews. For continuous characteristics, the regression coefficients obtained in the first stage were jointly pooled across reviews, using bivariate meta-analysis.^{19 20} For categorical characteristics the results of univariable meta-analyses are presented. We planned to perform multivariable analyses to assess the association between various study characteristics in combination and the performance of prognostic models, but due to the paucity of data we were not able to do so. All analyses are described in more detail in Supplement 3.

Patient and public involvement

Patients and public were not involved in the design, recruitment or conduct of the study.

Results

Identification and selection of studies

The search identified 2037 studies, of which 496 were included in the database and screened on full text, and 66 were further assessed (Figure 1). Finally, ten systematic reviews were included.²¹⁻²⁹ These reviews addressed external validations of the following prognostic models: ABCD2,³⁰ Essen Stroke Risk Score (ESRS),³¹ EuroSCORE,³² Framingham,³³ FRAX,³⁴ Injury Severity Score (ISS),³⁵ model for end-stage liver disease (MELD),³⁶ Pneumonia Severity Index (PSI),³⁷ Revised Cardiac Risk Index (RCRI),³⁸ and Simplified Acute Physiology Score (SAPS) 3³⁹ (Table 1). The reviews included 248 primary external validation studies with 274 external model validations (one study could describe multiple model validations). During data extraction, 73 of 274 validations were eventually excluded (most often for not reporting a performance measure), and 20 additional external model validations were identified (Figure 1). This resulted in the inclusion of 224 external validations, of which 221 could be included in the analyses of the c-statistic, and 124 in the analyses of the OE ratio. For the OE ratio, only validations of the EuroSCORE, Framingham, FRAX, PSI, RCRI and SAPS 3 prognostic models were included, due to the very low number of reported OE ratios in the validations studies for the other four prognostic models.

Description of included validations

The number of external validations per systematic review ranged from 11 to 30 (Table 1), and the median (IQR) sample size and number of events were 1069 (418-3043) and 92 (36-248), respectively. Most studies used an existing registry (N=104, 46%) or existing cohort (N=74, 33%) to validate the prognostic model. The median (IQR) c-statistic and OE ratio were 0.73 (0.64-0.82) and 0.92 (0.64-1.26), respectively. Predictive performance of the models was highly heterogeneous, even for external validations of the same prognostic model, as indicated by the wide prediction intervals (Table 1). Not all information on the study characteristics was reported for all external validations (Table S1). Information was often unclear (eg, for outcome definitions (N=83, 37%) and handling of missing data (N=105, 47%)) or missing (eg, case-mix information such as mean age (N=28, 13%) and gender distribution (N=16, 7%)).

Discrimination

Pooled models

1
2
3 250 The pooled analyses across all systematic reviews (Figure 2, S1 and S2) showed that validation in a
4
5 251 continent different from the development study was associated with a higher c-statistic, compared to
6
7 252 validation in the same continent, and multicentre versus single centre validation studies were associated
8
9 253 with a lower c-statistic. Comparable eligibility criteria for participant inclusion were also associated with
10
11 254 higher c-statistics compared to narrower criteria, whereas a broader setting was associated with a lower
12
13 255 c-statistic compared to a setting comparable to the development study. Although not statistically
14
15 256 significant, validations with changes made to the predictors (ie, substitution or deletion of a predictor),
16
17 257 or in which it was unclear whether all predictors were correctly measured, tended to have lower c-
18
19 258 statistics compared to validations where no changes were made. In various reviews we found an
20
21 259 association between the c-statistic and numerous other study characteristics, such as the study design,
22
23 260 comparability of outcome definition, prediction horizon, sample size and number of events, and mean
24
25 261 age of study participants (Figure 3, S2 and S3), only these were often not statistically significant when
26
27 262 pooled together.

28
29 263
30 264 *Variation across reviews*
31
32 265 Across reviews we found associations of many study characteristics with the c-statistic although this was
33
34 266 rather heterogeneous, and confidence intervals often overlapped (Figure 3 and Figure S3). For example,
35
36 267 for study design, in six systematic reviews a higher c-statistic was found for validations that used an
37
38 268 existing registry compared to an existing cohort, while in three reviews a lower c-statistic was found. In
39
40 269 three systematic reviews we found a higher c-statistic in validations by independent investigators, while
41
42 270 in five a lower c-statistic was found.

43
44 271
45 272 For other study characteristics, directions of associations were more consistent. For example, for most
46
47 273 systematic reviews, validation studies with eligibility criteria narrower compared to the criteria used in
48
49 274 the development study had a lower c-statistics while broader eligibility criteria were associated with
50
51 275 higher c-statistics (Figure S3). C-statistics were also (slightly) higher in external validations with a setting
52
53 276 comparable to the development study. Validation in a continent other than the development study in
54
55 277 general was associated with a higher c-statistic, and multicentre studies had lower c-statistics compared
56
57 278 to single centre studies. External validations in which it was unclear if there were changes made to the
58
59 279 predictors had lower c-statistics (Figure S3).

60 280
281 Calibration

282 *Pooled analyses*

283 We found a significant association between study design and the OE ratio (Figure 4); using data from a
284 case-control study (although known to be an inferior design for prognostic model research¹⁶) resulted in
285 higher OE ratios, compared to using data from an existing cohort (though based on three external
286 validations). Furthermore, higher OE ratios were found for studies in which the outcome was assessed
287 by a panel of clinicians as compared to using a registry. In various reviews we found an association
288 between the c-statistic and numerous other study characteristics, such as the duration of follow-up,
289 year in which recruitment was started, sample size, standard deviation of age, and setting (Figure 4, S4,
290 S5 and S6), only these were not statistically significant when pooled together.

292 *Variation across reviews*

293 For other categories of study design (other than the use of a case-control design), heterogeneous
294 associations were found across systematic reviews (Figure 5). The associations of most other study
295 characteristics with the OE ratio were also most often not consistent across systematic reviews (Figure
296 S5 and S6). For example, for two systematic reviews external validations with appropriate handling of
297 missing data had OE ratios closer to 1 compared to inappropriate handling of missing data, while in two
298 reviews, OE ratios were further away from 1. Only for the continent in which the model was validated,
299 directions were more consistent; OE ratios were closer to 1 if the continent was comparable to the
300 development, compared to validations in different continents (Figure S6).

Discussion

Principal findings

Using a comprehensive meta-analytical approach, we studied the association between study characteristics of prognostic model validation studies and the estimated model performance across ten clinical domains. We focused on objective study characteristics that can be extracted from published reports. The reporting of the primary external validation studies was often incomplete and inadequate. Key study characteristics, such as outcome definitions, handling of missing data, and even model calibration estimates were infrequently reported. Still, we found associations between various study characteristics and a model's predictive performance. Changes in a model's predictive performance were notably found in relation to validation studies with a case-control (versus cohort) design, with differences in case-mix, in continent (in which the model is validated), in eligibility criteria, in clinical setting, in number of centres (included in the validation study), in differences in outcome definitions and assessments, and in predictor substitutions.

Explanations, strengths and weaknesses

Based on findings in meta-epidemiological studies on the effect of study characteristics and the efficacy of interventions⁷⁻¹⁰ and diagnostic test accuracy,^{11 12} we anticipated to find more statistically significant associations between study characteristics and model performance across the included systematic reviews from different domains. Although we included every systematic review that described at least ten external validation studies of the same prognostic model, resulting in more than 200 validations from 10 reviews, our analyses appeared to still be hampered by relatively low numbers of external validations per systematic review, combined with poor reporting and substantial heterogeneity within and across systematic reviews.

Conceptually, there are many potential sources of heterogeneity in model performance, such as differences in population characteristics, predictor and outcome definitions and measurements, and in many aspects of the statistical analyses (eg, dealing with missing data, sample size and selective loss to follow up). All these characteristics may act in isolation but could also be related to each other. The individual strength of the association of one characteristic with model performance is ideally addressed by adopting multivariable (meta)-regression models with the observed model performance estimates of the validation studies as dependent variable and the characteristics of multiple design features as independent variables.^{10 12} Unfortunately, this approach was not feasible here due to the limited

number of available, well-reported, validation studies within the individual reviews, rendering the effective sample size too small for multivariable meta-regression analyses.

A general limitation of all meta-epidemiological studies, is the possibility that the effect of a certain study characteristic differs across systematic reviews which may nullify the effect when pooled together.⁴⁰ We also found numerous conflicting associations between a study characteristic and the reported predictive performance measures across reviews that were cancelled out in the pooled analyses.

Also, it is possible that the effect caused by individual study characteristics is small and therefore difficult to detect. Moreover, there might be some misclassification of study characteristics, caused either by our misinterpretation of what is reported, or by a lack of reporting, which could have diluted the effects of the study characteristics. Indeed, the c-statistic is often considered to be an insensitive measure to quantify changes in model performance.⁴¹⁻⁴³ In previous simulation studies, the c-statistic and OE ratio appeared to be strongly influenced by case-mix differences,^{14 44 45} which may mask the possible (smaller) effects from design-related characteristics. Other measures that are less sensitive to case-mix differences, such as the calibration slope, could, however, not be studied here simply because they were (almost) never reported in our retrieved studies, as was also shown previously.³

We found greater variation in the methods used by external validation studies *between* models than within validations of the *same* model. For example, multiple imputation is the preferred method for handling missing data in prediction modelling.^{46 47} However, in the field of cardiovascular disease, it seems common to handle missing data by performing a complete case analysis, while in the field of mortality prediction in surgical patients, typically researchers fill in 'normal' values if a value is missing. Finally, given the explorative nature of our analyses to identify potential areas of further research, we did not correct for multiple testing, though we tried to minimize the number of exploratory analyses.

Comparison to previous research

Despite above considerations, our findings, ie, the trends in the associations between study characteristics and model performance measures (though not always statistically significant), are in agreement with various previous simulation studies in this field.^{14 44 46-48} For example, we confirmed that studies with more variation in case-mix show higher c-statistics, and lower c-statistics when a predictor was omitted from the model. However, we found lower c-statistics in studies with a broader setting and when the number of centres in a study was higher.

1
2
3 364 We also found a higher OE ratio in studies with a case-control design. Both simulation studies and meta-
4 365 epidemiological studies in the fields of diagnostic tests and (mainly diagnostic) prediction models, have
5 366 shown biased effect measures in studies using a case-control design.¹¹⁻¹⁴ Further, we found that the OE
6 367 ratio was influenced by the method of outcome assessment, in agreement with previous studies that
7 368 showed that higher diagnostic odds ratios were found in studies with differential outcome verification.¹³
8 369 We also expected to find lower OE ratios when the validation population differed from the development
9 370 population (eg, in terms of case-mix).¹⁴ We could not systematically confirm this across all reviews, likely
10 371 caused by heterogeneity between systematic reviews as indicated by the wide confidence intervals.
11 372 Finally, we could not fully confirm the association between sample size and model performance that was
12 373 previously found,¹³ although we found similar trends in part of the reviews.
13
14
15
16
17
18
19
20 374

21 375 Implications for future research

22 376 In agreement with many previously conducted systematic reviews on prediction models,^{3 49-53} we still
23 377 and again found poor reporting of prediction model studies. Meta-epidemiological studies of prediction
24 378 model studies would highly benefit from complete reporting according to the Transparent reporting of a
25 379 multivariable prediction model for individual prognosis or diagnosis (TRIPOD) statement.^{16 17}
26 380 We also believe that more research is urgently needed to evaluate under which circumstances certain
27 381 design choices may lead to heterogeneity in prediction model performance, and to incorporate these
28 382 issues in the appraisal of prediction model studies. There is a need for more guidance on how to score
29 383 items of critical appraisal checklists for prediction model studies, such as the CHARMS checklist.¹⁵
30 384 Several options exist to gain more empirical insight in design related bias in prediction model studies.
31 385 Firstly, meta-epidemiological researchers can collect more external validation studies and try to correct
32 386 for all issues that cause variation in performance of a model. We believe, however, that this is currently
33 387 not feasible as we already included every systematic reviews describing at least ten validations of the
34 388 same prognostic model. A second and much more efficient option is to collect the individual participant
35 389 data (IPD) for all studies included in this review to directly study the effect of study characteristics on
36 390 model performance.⁵⁴⁻⁵⁸ Using IPD, it will also be possible to study different performance measures, like
37 391 the case-mix adjusted c-statistic^{44 59} and calibration slope.¹⁴ Thirdly, new simulation studies could be
38 392 performed to get more insight in design related bias in prediction model performance. Researchers
39 393 could for example study the effect of using a different outcome definition or prediction horizon on the c-
40 394 statistic of a model.
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55 395

Conclusion

In this comprehensive meta-epidemiological study we found empirical evidence for an association between study characteristics and predictive performance of prognostic models. We found that predictive performance of prognostic models upon external validation is highly heterogeneous, but sensitive to various study characteristics, such as study design, case-mix, eligibility criteria, setting, methods of outcome definition and measurement, and predictor substitution. It is important that these characteristics are thus emphasized in the reporting and appraisal of prediction model studies. However, for a large part the observed heterogeneity in model performance remained unexplained, which is likely caused by the high number of factors that cause heterogeneity in predictive performance and may act in opposite directions whereas a multivariable meta-regression analysis across reviews simply was not possible.

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

Acknowledgments: The authors would like to acknowledge Prof. Gary S Collins (GSC) for building the database with systematic reviews of prediction models, which served as a basis for this paper.

Contributors: KGMM, JBR, TPAD, and LH conceived the study. All authors were involved in designing the study. JAAGD selected the articles. JAAGD and RP extracted the data. JAAGD analysed the data in close consultation with TPAD. All authors were involved in interpreting the data. JAAGD wrote the first draft of the manuscript which was revised by all authors. The corresponding author attests that all listed authors meet authorship criteria and that no others meeting the criteria have been omitted. JAAGD is guarantor.

Funding: Thomas Debray gratefully acknowledges the Netherlands Organization for Health Research and Development (grant number 91617050). Karel GM Moons received a grant from The Netherlands Organization for Scientific Research (ZONMW 918.10.615 and 91208004). The funder had no role in the design of the study; the collection, analysis, and interpretation of the data; or approval of the finished manuscript.

Competing interests: All authors have completed the ICMJE uniform disclosure form at www.icmje.org/coi_disclosure.pdf and declare: no support from any organisation for the submitted work; no financial relationships with any organisations that might have an interest in the submitted work in the previous three years; no other relationships or activities that could appear to have influenced the submitted work.

Ethical approval: Not required.

Data sharing: no additional data are available.

Transparency: The lead author (JAAGD) affirms that the manuscript is an honest, accurate, and transparent account of the study being reported; that no important aspects of the study have been omitted; and that any discrepancies from the study as planned (and, if relevant, registered) have been explained.

1
2
3 438 The Corresponding Author has the right to grant on behalf of all authors and does grant on behalf of all
4
5 439 authors, a worldwide licence to the Publishers and its licensees in perpetuity, in all forms, formats and
6
7 440 media (whether known now or created in the future), to i) publish, reproduce, distribute, display and
8
9 441 store the Contribution, ii) translate the Contribution into other languages, create adaptations, reprints,
10
11 442 include within collections and create summaries, extracts and/or, abstracts of the Contribution, iii)
12
13 443 create any other derivative work(s) based on the Contribution, iv) to exploit all subsidiary rights in the
14
15 444 Contribution, v) the inclusion of electronic links from the Contribution to third party material where-
16
17 445 ever it may be located; and, vi) licence any third party to do any or all of the above.
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

References

1. Moons KG, Kengne AP, Woodward M, Royston P, Vergouwe Y, Altman DG, et al. Risk prediction models: I. Development, internal validation, and assessing the incremental value of a new (bio)marker. *Heart* 2012;98(9):683-90.

2. Debray TP, Vergouwe Y, Koffijberg H, Nieboer D, Steyerberg EW, Moons KG. A new framework to enhance the interpretation of external validation studies of clinical prediction models. *J Clin Epidemiol* 2014.

3. Damen JA, Hooft L, Schuit E, Debray TP, Collins GS, Tzoulaki I, et al. Prediction models for cardiovascular disease risk in the general population: systematic review. *BMJ* 2016;353:i2416.

4. Alba AC, Agoritsas T, Jankowski M, Courvoisier D, Walter SD, Guyatt GH, et al. Risk prediction models for mortality in ambulatory patients with heart failure: a systematic review. *Circ Heart Fail* 2013;6(5):881-9.

5. Meads C, Ahmed I, Riley RD. A systematic review of breast cancer incidence risk prediction models with meta-analysis of their performance. *Breast Cancer Res Treat* 2012;132(2):365-77.

6. Moons KG, Royston P, Vergouwe Y, Grobbee DE, Altman DG. Prognosis and prognostic research: what, why, and how? *BMJ* 2009;338:b375.

7. Page MJ, Higgins JP, Clayton G, Sterne JA, Hrobjartsson A, Savovic J. Empirical Evidence of Study Design Biases in Randomized Trials: Systematic Review of Meta-Epidemiological Studies. *PLoS One* 2016;11(7):e0159267.

8. Berkman ND, Santaguida PL, Viswanathan M, Morton SC. AHRQ Methods for Effective Health Care. *The Empirical Evidence of Bias in Trials Measuring Treatment Differences*. Rockville (MD): Agency for Healthcare Research and Quality (US), 2014.

9. Savovic J, Jones H, Altman D, Harris R, Juni P, Pildal J, et al. Influence of reported study design characteristics on intervention effect estimates from randomised controlled trials: combined analysis of meta-epidemiological studies. *Health Technol Assess* 2012;16(35):1-82.

10. Wood L, Egger M, Gluud LL, Schulz KF, Juni P, Altman DG, et al. Empirical evidence of bias in treatment effect estimates in controlled trials with different interventions and outcomes: meta-epidemiological study. *BMJ* 2008;336(7644):601-5.

11. Lijmer JG, Mol BW, Heisterkamp S, Bossel GJ, Prins MH, van der Meulen JH, et al. Empirical evidence of design-related bias in studies of diagnostic tests. *JAMA* 1999;282(11):1061-6.

12. Rutjes AW, Reitsma JB, Di Nisio M, Smidt N, van Rijn JC, Bossuyt PM. Evidence of bias and variation in diagnostic accuracy studies. *CMAJ* 2006;174(4):469-76.

13. Ban JW, Emparanza JI, Urreta I, Burls A. Design Characteristics Influence Performance of Clinical Prediction Rules in Validation: A Meta-Epidemiological Study. *PLoS One* 2016;11(1):e0145779.

14. Steyerberg E. *Clinical prediction models: a practical approach to development, validation, and updating*: Springer Science & Business Media, 2008.

15. Moons KG, de Groot JA, Bouwmeester W, Vergouwe Y, Mallett S, Altman DG, et al. Critical appraisal and data extraction for systematic reviews of prediction modelling studies: the CHARMS checklist. *PLoS Med* 2014;11(10):e1001744.

16. Collins GS, Reitsma JB, Altman DG, Moons KG. Transparent Reporting of a multivariable prediction model for Individual Prognosis or Diagnosis (TRIPOD): the TRIPOD statement. *Ann Intern Med* 2015;162(1):55-63.

17. Moons KG, Altman DG, Reitsma JB, Ioannidis JP, 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. Snell KI, Ensor J, Debray TP, Moons KG, Riley RD. Meta-analysis of prediction model performance across multiple studies: Which scale helps ensure between-study normality for the C-statistic and calibration measures? *Stat Methods Med Res* 2017;962280217705678.
19. Snell KI, Hua H, Debray TP, Ensor J, Look MP, Moons KG, et al. Multivariate meta-analysis of individual participant data helped externally validate the performance and implementation of a prediction model. *J Clin Epidemiol* 2015.
20. Debray TP, Damen JA, Snell KI, Ensor J, Hooft L, Reitsma JB, et al. A guide to systematic review and meta-analysis of prediction model performance. *BMJ* 2017;356:i6460.
21. Thompson DD, Murray GD, Dennis M, Sudlow CL, Whiteley WN. Formal and informal prediction of recurrent stroke and myocardial infarction after stroke: a systematic review and evaluation of clinical prediction models in a new cohort. *BMC Med* 2014;12:58.
22. Siregar S, Groenwold RH, de Heer F, Bots ML, van der Graaf Y, van Herwerden LA. Performance of the original EuroSCORE. *Eur J Cardiothorac Surg* 2012;41(4):746-54.
23. Marques A, Ferreira RJ, Santos E, Loza E, Carmona L, da Silva JA. The accuracy of osteoporotic fracture risk prediction tools: a systematic review and meta-analysis. *Ann Rheum Dis* 2015;74(11):1958-67.
24. Tohira H, Jacobs I, Mountain D, Gibson N, Yeo A. Systematic review of predictive performance of injury severity scoring tools. *Scand J Trauma Resusc Emerg Med* 2012;20:63.
25. Klein KB, Stafinski TD, Menon D. Predicting survival after liver transplantation based on pre-transplant MELD score: a systematic review of the literature. *PLoS One* 2013;8(12):e80661.
26. Chalmers JD, Mandal P, Singanayagam A, Akram AR, Choudhury G, Short PM, et al. Severity assessment tools to guide ICU admission in community-acquired pneumonia: systematic review and meta-analysis. *Intensive Care Med* 2011;37(9):1409-20.
27. Ford MK, Beattie WS, Wijeyesundera DN. Systematic review: prediction of perioperative cardiac complications and mortality by the revised cardiac risk index. *Ann Intern Med* 2010;152(1):26-35.
28. Nassar AP, Malbouisson LM, Moreno R. Evaluation of Simplified Acute Physiology Score 3 performance: a systematic review of external validation studies. *Crit Care* 2014;18(3):R117.
29. Damen JAAG, Pajouheshnia R, Heus P, Moons KGM, Reitsma JB, Scholten RJPM, et al. Performance of the Framingham risk models and Pooled Cohort Equations for predicting 10-year risk of cardiovascular disease: a systematic review and meta-analysis. Manuscript submitted for publication.
30. Rothwell PM, Giles MF, Flossmann E, Lovelock CE, Redgrave JN, Warlow CP, et al. A simple score (ABCD) to identify individuals at high early risk of stroke after transient ischaemic attack. *Lancet* 2005;366(9479):29-36.
31. Diener HC, Ringleb PA, Savi P. Clopidogrel for the secondary prevention of stroke. *Expert Opin Pharmacother* 2005;6(5):755-64.
32. Nashef SA, Roques F, Michel P, Gauducheau E, Lemeshow S, Salamon R. European system for cardiac operative risk evaluation (EuroSCORE). *Eur J Cardiothorac Surg* 1999;16(1):9-13.
33. Wilson PW, D'Agostino RB, Levy D, Belanger AM, Silbershatz H, Kannel WB. Prediction of coronary heart disease using risk factor categories. *Circulation* 1998;97(18):1837-47.
34. Kanis JA, Oden A, Johnell O, Johansson H, De Laet C, Brown J, et al. The use of clinical risk factors enhances the performance of BMD in the prediction of hip and osteoporotic fractures in men and women. *Osteoporos Int* 2007;18(8):1033-46.
35. Baker SP, O'Neill B, Haddon W, Jr., Long WB. The injury severity score: a method for describing patients with multiple injuries and evaluating emergency care. *J Trauma* 1974;14(3):187-96.

1
2
3 538 36. Malinchoc M, Kamath PS, Gordon FD, Peine CJ, Rank J, ter Borg PC. A model to predict poor survival
4 539 in patients undergoing transjugular intrahepatic portosystemic shunts. *Hepatology*
5 540 2000;31(4):864-71.
6 541 37. Fine MJ, Auble TE, Yealy DM, Hanusa BH, Weissfeld LA, Singer DE, et al. A prediction rule to identify
7 542 low-risk patients with community-acquired pneumonia. *N Engl J Med* 1997;336(4):243-50.
8 543 38. Lee TH, Marcantonio ER, Mangione CM, Thomas EJ, Polanczyk CA, Cook EF, et al. Derivation and
9 544 prospective validation of a simple index for prediction of cardiac risk of major noncardiac
10 545 surgery. *Circulation* 1999;100(10):1043-9.
11 546 39. Moreno RP, Metnitz PG, Almeida E, Jordan B, Bauer P, Campos RA, et al. SAPS 3--From evaluation of
12 547 the patient to evaluation of the intensive care unit. Part 2: Development of a prognostic model
13 548 for hospital mortality at ICU admission. *Intensive Care Med* 2005;31(10):1345-55.
14 549 40. Sterne JA, Juni P, Schulz KF, Altman DG, Bartlett C, Egger M. Statistical methods for assessing the
15 550 influence of study characteristics on treatment effects in 'meta-epidemiological' research. *Stat*
16 551 *Med* 2002;21(11):1513-24.
17 552 41. Steyerberg EW, Vickers AJ, Cook NR, Gerds T, Gonen M, Obuchowski N, et al. Assessing the
18 553 performance of prediction models: a framework for traditional and novel measures.
19 554 *Epidemiology* 2010;21(1):128-38.
20 555 42. Pencina MJ, D'Agostino RB, Sr., D'Agostino RB, Jr., Vasan RS. Evaluating the added predictive ability
21 556 of a new marker: from area under the ROC curve to reclassification and beyond. *Stat Med*
22 557 2008;27(2):157-72; discussion 207-12.
23 558 43. Uno H, Cai T, Pencina MJ, D'Agostino RB, Wei LJ. On the C-statistics for evaluating overall adequacy
24 559 of risk prediction procedures with censored survival data. *Stat Med* 2011;30(10):1105-17.
25 560 44. Vergouwe Y, Moons KG, Steyerberg EW. External validity of risk models: Use of benchmark values to
26 561 disentangle a case-mix effect from incorrect coefficients. *Am J Epidemiol* 2010;172(8):971-80.
27 562 45. Usher-Smith JA, Sharp SJ, Griffin SJ. The spectrum effect in tests for risk prediction, screening, and
28 563 diagnosis. *BMJ* 2016;353:i3139.
29 564 46. Held U, Kessels A, Garcia Aymerich J, Basagana X, Ter Riet G, Moons KG, et al. Methods for Handling
30 565 Missing Variables in Risk Prediction Models. *Am J Epidemiol* 2016;184(7):545-51.
31 566 47. Janssen KJ, Vergouwe Y, Donders AR, Harrell FE, Jr., Chen Q, Grobbee DE, et al. Dealing with missing
32 567 predictor values when applying clinical prediction models. *Clin Chem* 2009;55(5):994-1001.
33 568 48. Austin PC, Steyerberg EW. Interpreting the concordance statistic of a logistic regression model:
34 569 relation to the variance and odds ratio of a continuous explanatory variable. *BMC Med Res*
35 570 *Methodol* 2012;12:82.
36 571 49. Collins GS, Mallett S, Omar O, Yu LM. Developing risk prediction models for type 2 diabetes: a
37 572 systematic review of methodology and reporting. *BMC Med* 2011;9:103.
38 573 50. Collins GS, Omar O, Shanyinde M, Yu LM. A systematic review finds prediction models for chronic
39 574 kidney disease were poorly reported and often developed using inappropriate methods. *J Clin*
40 575 *Epidemiol* 2013;66(3):268-77.
41 576 51. Bouwmeester W, Zuithoff NP, Mallett S, Geerlings MI, Vergouwe Y, Steyerberg EW, et al. Reporting
42 577 and methods in clinical prediction research: a systematic review. *PLoS Med* 2012;9(5):1-12.
43 578 52. Perel P, Edwards P, Wentz R, Roberts I. Systematic review of prognostic models in traumatic brain
44 579 injury. *BMC Med Inform Decis Mak* 2006;6:38.
45 580 53. Mallett S, Royston P, Dutton S, Waters R, Altman DG. Reporting methods in studies developing
46 581 prognostic models in cancer: a review. *BMC Med* 2010;8:20.
47 582 54. Debray TP, Koffijberg H, Vergouwe Y, Moons KG, Steyerberg EW. Aggregating published prediction
48 583 models with individual participant data: a comparison of different approaches. *Stat Med*
49 584 2012;31(23):2697-712.

55. Debray TP, Moons KG, Ahmed I, Koffijberg H, Riley RD. A framework for developing, implementing, and evaluating clinical prediction models in an individual participant data meta-analysis. *Stat Med* 2013;32(18):3158-80.
56. Debray TP, Riley RD, Rovers MM, Reitsma JB, Moons KG. Individual Participant Data (IPD) Meta-analyses of Diagnostic and Prognostic Modeling Studies: Guidance on Their Use. *PLoS Med* 2015;12(10):e1001886.
57. Riley RD, Ensor J, Snell KI, Debray TP, Altman DG, Moons KG, et al. External validation of clinical prediction models using big datasets from e-health records or IPD meta-analysis: opportunities and challenges. *BMJ* 2016;353:i3140.
58. Riley RD, Price MJ, Jackson D, Wardle M, Gueyffier F, Wang J, et al. Multivariate meta-analysis using individual participant data. *Res Synth Methods* 2015;6(2):157-74.
59. White IR, Rapsomaniki E. Covariate-adjusted measures of discrimination for survival data. *Biom J* 2015;57(4):592-613.
60. Giles MF, Rothwell PM. Systematic review and pooled analysis of published and unpublished validations of the ABCD and ABCD2 transient ischemic attack risk scores. *Stroke* 2010;41(4):667-73.

Table 1: Description of included reviews and prediction models

Systematic review	Giles 2010 ⁶⁰	Thompson 2014 ²¹	Siregar 2012 ²²	Damen ²⁹	Marques 2015 ²³	Tohira 2012 ²⁴	Klein 2013 ²⁵	Chalmers 2011 ²⁶	Ford 2010 ²⁷	Nassar 2014 ²⁸
Model	ABCD2 ³⁰	ESRS ³¹	EuroSCORE ³²	Framingham ³	FRAX ³⁴	ISS ³⁵	MELD ³⁶	PSI ³⁷	RCRI ³⁸	SAPS 3 ³⁹
Population	Patients with TIA	Adults with a previous CVD event	Adult patients who underwent cardiac surgery under cardiopulmonary bypass	Men without previous CHD event	General population	Injured patients	Patients with liver cirrhosis but without hepatocellular carcinoma who underwent elective transjugular intrahepatic portosystemic shunts	Inpatients with community-acquired pneumonia	Patients aged >=50 years who underwent nonemergent noncardiac procedures	ICU patients
Geographical location	United States and UK	Canada, United States, Europe	Europe	United States	Europe, Canada, Japan, United States, Australia	United States	United States	United States	United States	Worldwide
Patient recruitment	1981-1998	1992-1995	1995	1971-1974	1980-1999	1968-1969	1991-1995	1989	1989-1994	2002
Predicted outcome	Stroke	Recurrent ischemic stroke, MI and vascular death	Mortality	CHD	Osteoporotic fractures	All-cause mortality	All-cause mortality	30-day hospital mortality	Major cardiac complications	Hospital mortality
Prediction horizon	2 days	1 year	30 days	10 years	10 years	3 months	3 months	30 days	1 year	90 days
Performance development study										
C-statistic	0.66 [95% CI 0.60, 0.71]	NR	0.7875	0.74	0.63	NR	NR	0.84	0.759 [SE 0.032]	0.848
OE ratio	NR*	NR*	NR*	NR*	NR*	NR*	NR*	NR*	NR*	1.00 [95% CI 0.98, 1.02]
Pooled performance validation studies										
Number of external validations included in analyses	16	11	22	23	30	34	14	24	23	27
C-statistic [95% CI]	0.66 [0.61, 0.71]	0.60 [0.58, 0.62]	0.79 [0.77, 0.81]	0.68 [0.65, 0.71]	0.66 [0.63, 0.68]	0.86 [0.83, 0.88]	0.64 [0.59, 0.68]	0.80 [0.77, 0.82]	0.69 [0.65, 0.72]	0.83 [0.80, 0.85]
95% PI	[0.54, 0.77]	[0.57, 0.63]	[0.74, 0.83]	[0.56, 0.78]	[0.54, 0.76]	[0.62, 0.96]	[0.48, 0.77]	[0.64, 0.89]	[0.53, 0.81]	[0.66, 0.92]
OE ratio [95% CI]	NA	NA	0.54 [0.42, 0.68]	0.58 [0.45, 0.76]	1.10 [0.83, 1.47]	NA	NA	0.94 [0.83, 1.06]	2.70 [1.72, 4.25]	0.89 [0.77, 1.03]

95% PI	NA	NA	[0.19, 1.51]	[0.20, 1.74]	[0.31, 3.93]	NA	NA	[0.55, 1.60]	[0.35, 20.75]	[0.42, 1.91]
--------	----	----	--------------	--------------	--------------	----	----	--------------	---------------	--------------

TIA: transient ischaemic attack, CVD: cardiovascular disease, CHD: coronary heart disease, ICU: intensive care unit, UK: United Kingdom, MI: myocardial infarction, NR: not reported, CI: confidence interval, PI: prediction interval, NA: not assessed.

*As the models are optimally fit in the development dataset, all OE ratios should be close to 1.

For peer review only

Figure legends

Figure 1: Flow chart of study selection.
SR: systematic review, IPD: individual participant data, MA: meta-analysis, NR: not reported, c: concordance, OE: observed expected.

Figure 2: Associations between study characteristics and logit c-statistic with regard to a reference category across 221 external validation studies and 10 different prediction models. Figure S1 shows these differences on the original scale if we assume a c-statistic of 0.70 in the reference category. For example, for comparability of eligibility criteria, if we assume a c-statistic of 0.70 in the reference category (narrower), this would result in c-statistics of 0.74 [0.72, 0.77], 0.73 [0.66, 0.79], 0.77 [0.68, 0.84], and 0.77 [0.59,0.89] in the categories comparable, mixture, broader, and unclear, respectively.

Figure 3: C-statistic for categories of study design, pooled using univariable meta-regression analyses within each systematic review. N represents the number of external validation studies in a specific category. C diff represents the difference in c-statistic with regard to a reference category (indicated with 'ref').

Figure 4: Associations between study characteristics and ln OE ratio with regard to a reference category across 124 external validation studies and 6 different prediction models. Figure S4 shows these differences on the original scale if we assume an OE ratio of 1.00 in the reference category. For example, for comparability of eligibility criteria, if we assume an OE ratio of 1.00 in the reference category (narrower), this would result in OE ratios of 0.83 [0.66, 1.05], 1.11 [0.72, 1.70], and 0.92 [0.54, 1.58] in the categories comparable, mixture, and broader, respectively.

Figure 5: OE ratio for categories of study design, pooled using univariable meta-regression analyses within each systematic review. N represents the number of external validation studies in a specific category. OE diff represents the difference in OE ratio with regard to a reference category (indicated with 'ref').

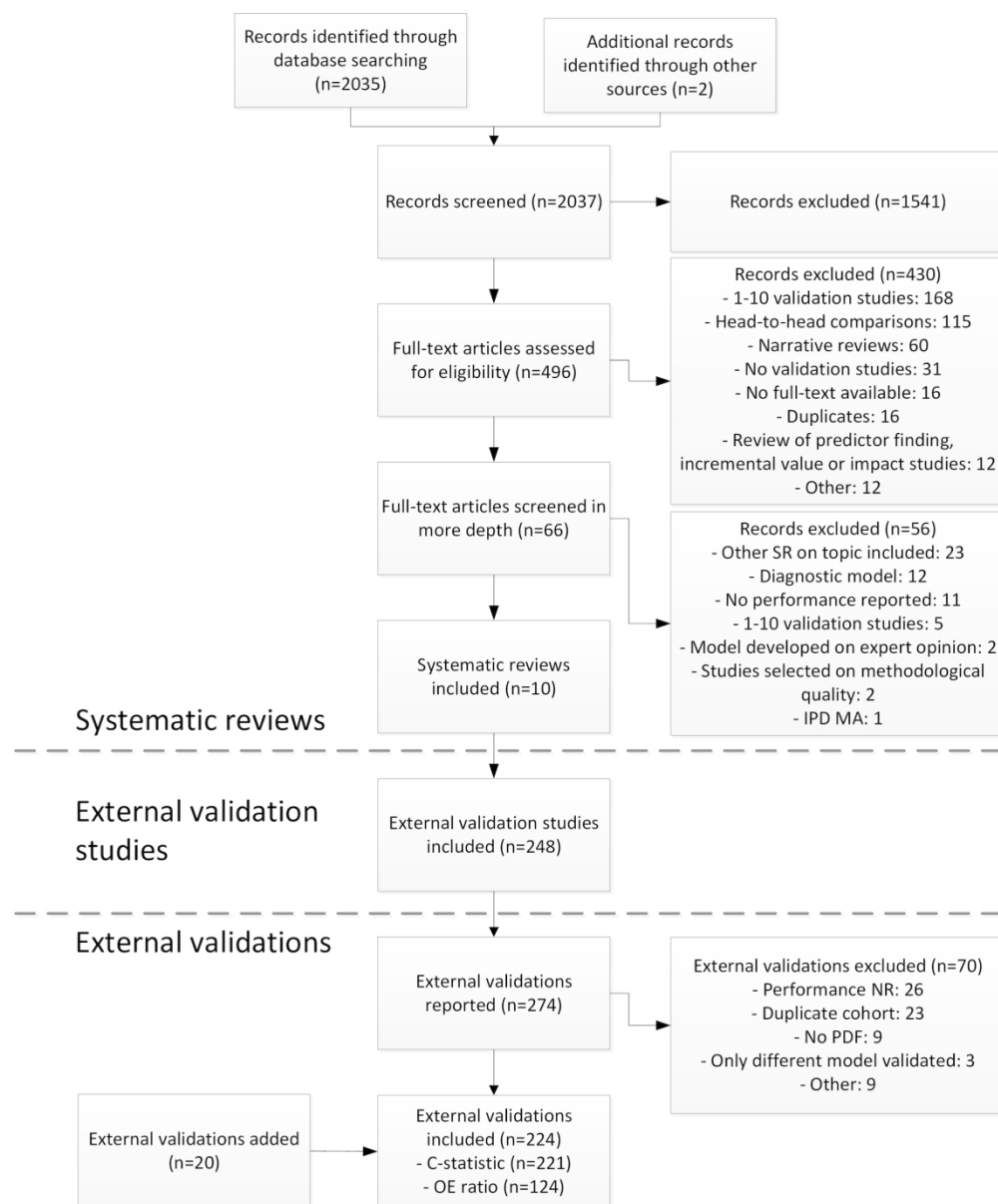


Figure 1: Flow chart of study selection.

SR: systematic review, IPD: individual participant data, MA: meta-analysis, NR: not reported, c: concordance, OE: observed expected.

184x222mm (300 x 300 DPI)

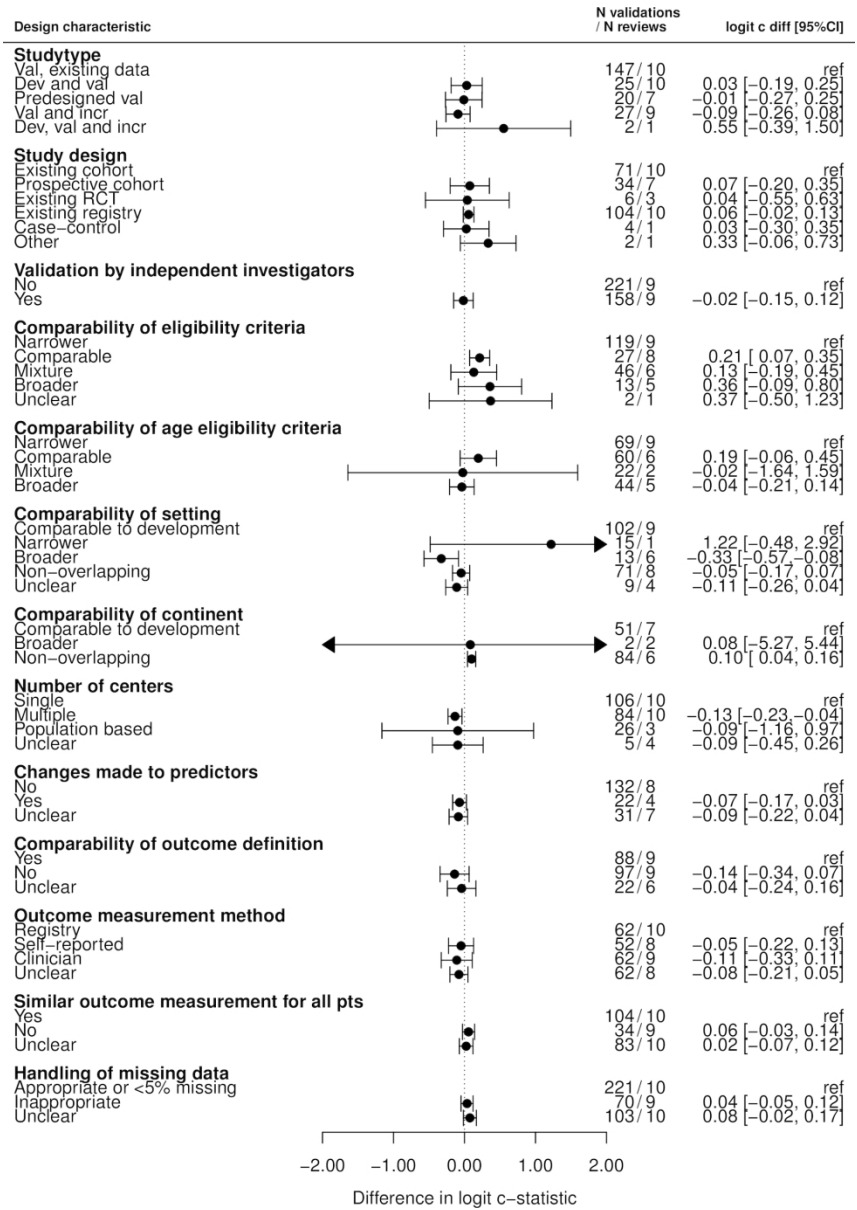


Figure 2: Associations between study characteristics and logit c-statistic with regard to a reference category across 221 external validation studies and 10 different prediction models. Figure S1 shows these differences on the original scale if we assume a c-statistic of 0.70 in the reference category. For example, for comparability of eligibility criteria, if we assume a c-statistic of 0.70 in the reference category (narrower), this would result in c-statistics of 0.74 [0.72, 0.77], 0.73 [0.66, 0.79], 0.77 [0.68, 0.84], and 0.77 [0.59, 0.89] in the categories comparable, mixture, broader, and unclear, respectively.

157x222mm (300 x 300 DPI)

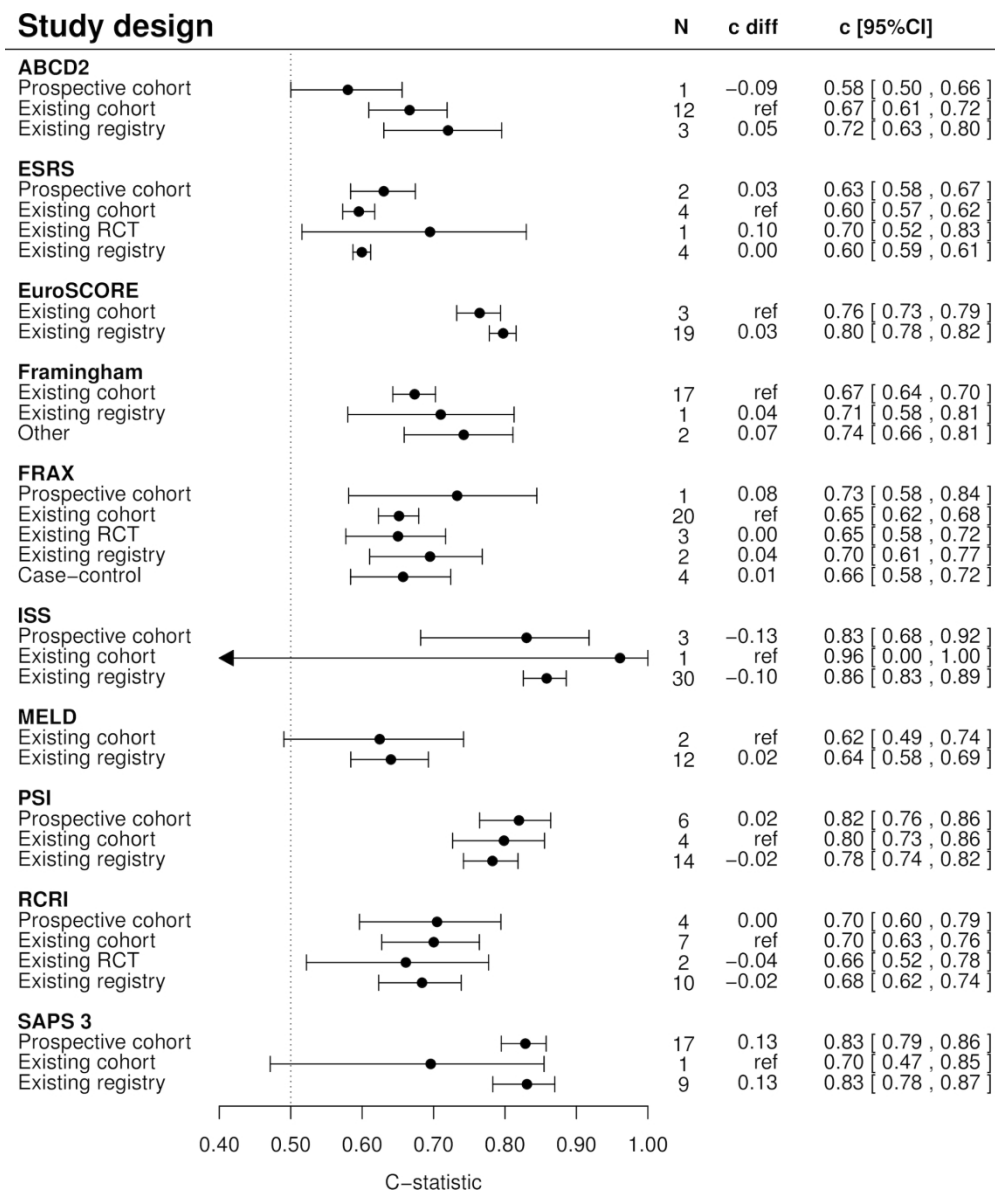


Figure 3: C-statistic for categories of study design, pooled using univariable meta-regression analyses within each systematic review. N represents the number of external validation studies in a specific category. C diff represents the difference in c-statistic with regard to a reference category (indicated with 'ref').

188x222mm (300 x 300 DPI)

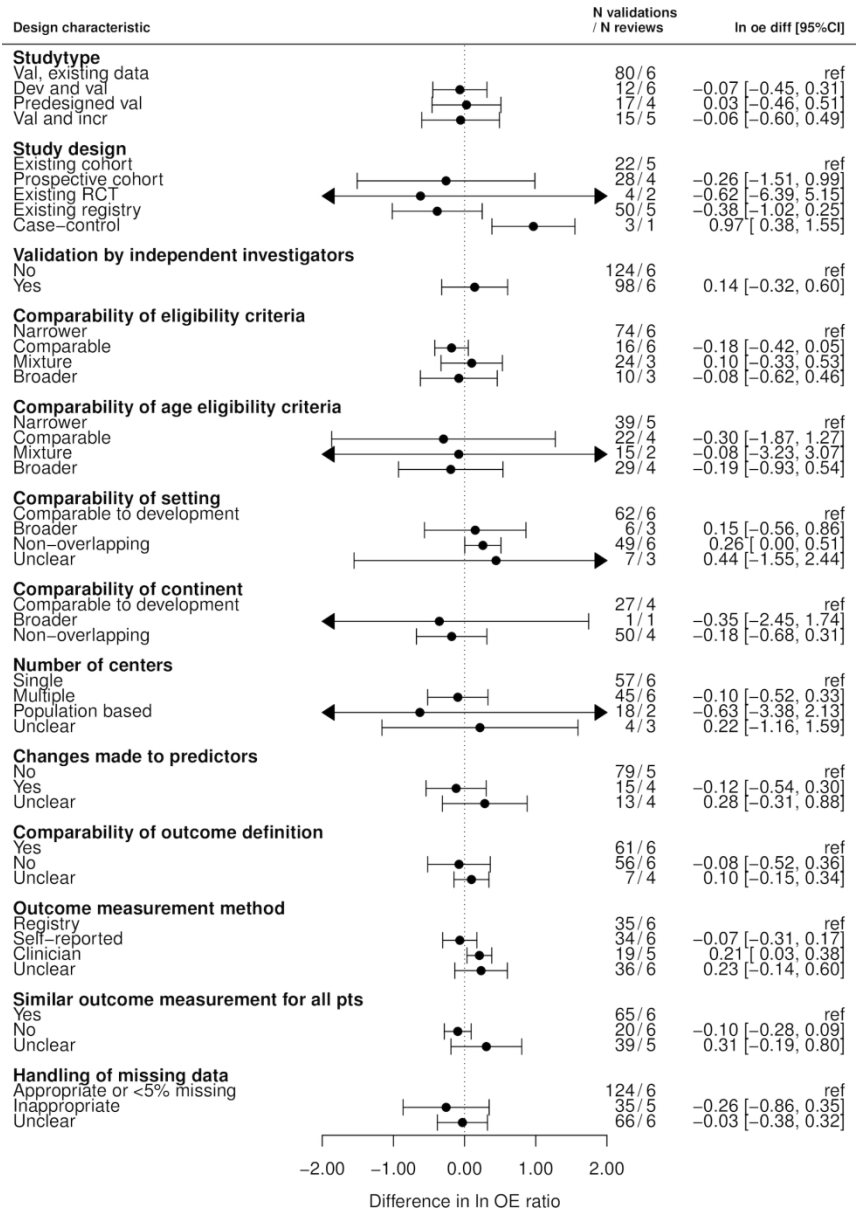


Figure 4: Associations between study characteristics and In OE ratio with regard to a reference category across 124 external validation studies and 6 different prediction models. Figure S4 shows these differences on the original scale if we assume an OE ratio of 1.00 in the reference category. For example, for comparability of eligibility criteria, if we assume an OE ratio of 1.00 in the reference category (narrower), this would result in OE ratios of 0.83 [0.66, 1.05], 1.11 [0.72, 1.70], and 0.92 [0.54, 1.58] in the categories comparable, mixture, and broader, respectively.

158x222mm (300 x 300 DPI)

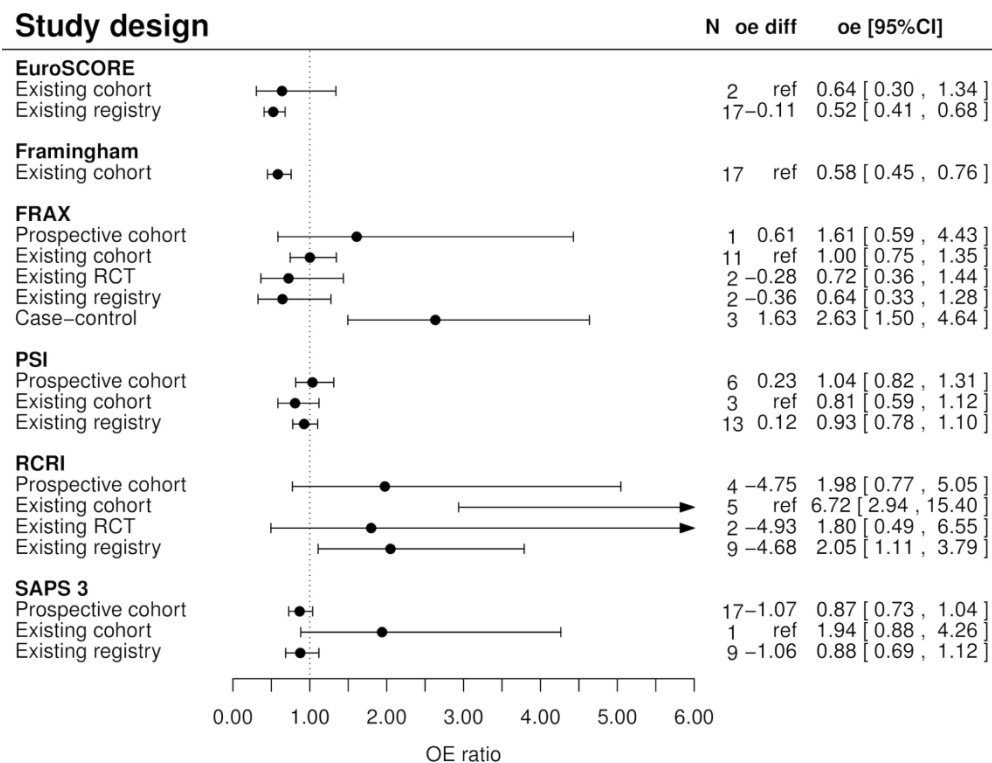


Figure 5: OE ratio for categories of study design, pooled using univariable meta-regression analyses within each systematic review. N represents the number of external validation studies in a specific category. OE diff represents the difference in OE ratio with regard to a reference category (indicated with 'ref').

190x142mm (300 x 300 DPI)

1

2

3 **Content**

4

5 Methods

- 6 Supplement 1: Search string
- 7 Supplement 2: Description of items extracted from studies and included in analyses
- 8 Supplement 3: Statistical analyses
- 9

10

11 Tables and figures

- 12 Table S1: Description of study characteristics and quality of reporting within each systematic review
- 13 Figure S1: Associations between categorical variables and c-statistic
- 14 Figure S2: Associations between continuous variables and c-statistic
- 15 Figure S3: C-statistic in categories of study characteristics within each systematic review
- 16 Figure S4: Associations between categorical variables and OE ratio
- 17 Figure S5: Associations between continuous variables and OE ratio
- 18 Figure S6: OE ratio in categories of study characteristics within each systematic review
- 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

Supplement 1: Search string

(clinical prediction[ti] OR
risk calculator*[ti] OR
risk index[ti] OR
risk indices[ti] OR
risk model*[ti] OR
risk prediction[ti] OR
risk score*[ti] OR
risk stratification[ti] OR
predictive model*[ti] OR
prediction model*[ti] OR
prediction rule*[tiab] OR
prognostic index[ti] OR
prognostic indices[ti] OR
prognostic model*[ti] OR
scoring system*[ti]) AND
(review[Publication Type] OR
review[ti] OR
critical appraisal[ti] OR
Bibliography[Publication Type] OR
Meta-analysis[Publication Type]) NOT
(Editorial[Publication Type] OR
Letter[Publication Type] OR
News[Publication Type])

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

Supplement 2: Description of items extracted from studies and included in analyses

Item	Extracted from studies	Categorization / handling in analyses	Description / examples
Validated model	ABCD2, ESRS, EuroSCORE, Framingham Wilson, FRAX, ISS, MELD, PSI, RCRI, SAPS 3	-	-
Study type	Predesigned validation study	Predesigned validation study	Study designed with the aim of validating a prediction model
	Validation study using existing data	Validation study using existing data	Study in which a prediction model is validated using a dataset collected for a different purpose than validating the model
	Development of new model and validation of different model	Development of new model and validation of different model	Study in which a model is developed and a model is validated
	Validation and incremental value	Validation and incremental value	Study in which a model is validated and in which the added value of one or more predictors is assessed
	Development, validation, and incremental value study	Development, validation, and incremental value study	Combination of the two above
Independent investigators	Yes	Yes	None of the authors of the development study was listed as author in the external validation study
	No	No	One or more of the authors of the development study was listed as author in the external validation study
Study design	Prospective cohort	Prospective cohort	
	Existing cohort	Existing cohort	
	Existing RCT	Existing RCT	
	Existing registry / medical records	Existing registry	
	Case-control	Case-control	
	Other (specify)	Other	
Eligibility criteria for participants	Copy/paste eligibility criteria of validation study	Comparable	Eligibility criteria comparable to development study

Item	Extracted from studies	Categorization / handling in analyses	Description / examples
		Narrower	People included in the development study excluded in the validation study
		Broader	People excluded in the development study included in the validation study
		Mixture	Combination of narrower and broader
		Unclear	
Setting	Primary care Secondary care Tertiary care Population based Screening Mixed Unclear	Comparable	Same setting as development study
		Broader	Same setting as development study, and participants from additional settings recruited
		Non-overlapping	Setting in development study differs from validation study
		Unclear	
Study dates	Start year of recruitment End year of recruitment	Continuous, standardized per systematic review	
Prediction horizon	Time period for which predictions were made, eg, 10 years.	Continuous, standardized per systematic review	
Geographical location	Country and continent	Comparable	Model validated in the same continent as the development study
		Broader	Model validated in the same and additional continents as the development study
		Non-overlapping	Model validated in a different continent than the development study
Number of centres	Number of centres (numerical)	Single	
		Multiple	
		Population based	Participants not recruited at medical centres, but, for example, from a specific geographic area (eg, all individuals living in Framingham, US)
		Unclear	
Case-mix: age mean	Mean and SD of age of	Continuous, standardized per	

Item	Extracted from studies	Categorization / handling in analyses	Description / examples
and sd	participants included in the study, or other available information about age distribution	systematic review	
Case-mix: gender	Percentage of men included in a study	Continuous, standardized per systematic review	
Predictors	Were predictors deleted from the model, or were predictors substituted with different predictors.	Yes	Changes made to predictors
		No	No changes made to predictors
		Unclear	
Predicted outcome	Full definition, including ICD-codes	Comparable	Outcome definition comparable to development study
		Not comparable	Outcome definition not comparable to development study
		Unclear	
Outcome - measurement method	Measurement method (eg, self-reported, interviews, expert panel), differences in outcome measurement between participants in the study	Yes	Outcome measurement similar for all participant
		No	Systematic differences in outcome measurement between participants
		Unclear	
Missing data	Number of participants with missing data, method of handling missing data	Appropriate	Missing data handled using multiple imputation, or <5% missing data (arbitrary cut-off)
		Inappropriate	Missing data not handled using multiple imputation (eg, complete-case analysis, mean imputation), and >=5% missing data
		Unclear	Unclear handling of missing data, and >=5% missing data
Number of participants		Continuous, standardized per systematic review	
Number of events		Continuous, standardized per systematic review	

Item	Extracted from studies	Categorization / handling in analyses	Description / examples
Model updating	Was the model altered before validating it, eg, using intercept recalibration.	NA	
Performance - c-statistic	C-statistic, AUC, 95% confidence intervals or SE	Logit transformation ¹	
Performance - OE ratio	OE ratio, predicted risks, presence of calibration plots or tables, 95% confidence intervals or SE	Ln transformation ¹	

SD: standard deviation, NA: not applicable, C-statistic: concordance statistic, AUC: area under the receiver operating curve, SE: standard error, OE ratio: observed expected ratio.

Information regarding c-statistics and OE ratios when not reported was sometimes restored from other information reported in the paper. If the precision of the c-statistic was not reported, we estimated this from the c-statistic and sample size of the study, using the formula described by Newcombe and Hanley.^{2,3} Various equations were used to estimate the standard error of the OE ratio, depending on which information was reported. All equations (as numbered) are described in the appendix of Debray et al.⁴ If the SE of the OE ratio was reported, we used equation 16 to estimate the SE of $\ln(\text{OE})$, if the observed event risk (P_o), the expected event risk (P_e), and the SE of P_o were reported, we used equation 51, and if only P_o and P_e were reported we used equation 27.

Supplement 3: Statistical analyses

First we pooled the total OE ratio and c-statistic within each systematic review. Based on previous recommendations,¹⁴ we pooled the log OE ratio and logit c-statistic using random-effects meta-analysis accounting for the presence of between-study heterogeneity, weighted by the inverse of the variance. We calculated 95% confidence intervals (CI) and (approximate) 95% prediction intervals (PI) to quantify uncertainty and the presence of between-study heterogeneity. The Hartung-Knapp-Sidik-Jonkman (HKSJ) method was used when calculating 95% CIs.⁵ The 95% PI was calculated using the equation described previously.⁴ The CI indicates the precision of the summary performance estimate and the PI provides boundaries on the likely performance in future model validation studies that are comparable to the studies included in the meta-analysis, and can thus be seen as an indication of model generalizability.⁶

To study the possible association between study characteristics and predictive performance, we used a two-stepped approach. In the first stage, we fitted a univariable meta-regression model (ie, a separate model for every study characteristic) within every systematic review, with the logit c-statistic or log OE ratio as outcome variable. This model was fitted with intercept term. Therefore, the effect estimates obtained from this meta-regression model indicate the difference in logit c-statistic or log OE ratio between a certain category of a study characteristic and a chosen reference category of that characteristic. As a reference category, we chose the category that was present in the highest number of systematic reviews allowing the inclusion of as many data as possible.

In the second stage, these effect estimates were pooled with a random effects meta-analysis model. This reflected the influence of the study characteristic on model performance over all systematic reviews. For continuous study characteristics, the intercept term and beta-coefficient from the first stage were jointly pooled across reviews using bivariate meta-analysis.^{4 6} For categorical study characteristics the data available were not sufficient for the complexity of a multivariate model, so every category was pooled in a separate (univariate) meta-analysis.

As the estimates obtained with this approach are on the transformed scale (ie, the difference in logit c-statistic or log OE ratio between one category and the reference category), we transformed these back assuming a c-statistic of 0.7 or an OE ratio of 1.00 in the reference category. Also, we performed a second analysis where we again fitted a univariable meta-regression model, with the logit c-statistic or log OE ratio as outcome variable, but now without intercept term. This analysis enables the calculation

1
2
3 of an effect estimate for every category of a study characteristic and to back transform this to the
4 original scale, yielding a pooled c-statistic or pooled OE ratio for each category of a study characteristic.
5
6 We planned to perform multivariable analyses to assess the association between various study
7 characteristics in combination and the performance of prediction models, but due to the paucity of data
8 we were not able to do so.
9
10 All analyses were performed in R version 3.3.2,⁷ using the packages metafor,⁸ mvmeta,⁹ metamisc,¹⁰ and
11 lme4.¹¹
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

Table S1: Description of study characteristics and quality of reporting within each systematic review

Categorical variables

	ABCD2	ESRS	EuroSCORE	Framingham	FRAX	ISS	MELD	PSI	RCRI	SAPS 3
Studytype										
Validation study using existing data	9 (56%)	7 (64%)	21 (95%)	18 (78%)	26 (87%)	24 (71%)	10 (71%)	16 (67%)	11 (48%)	8 (30%)
Development of new model and validation of different model	2 (12%)	2 (18%)	1 (5%)	2 (9%)	1 (3%)	4 (12%)	3 (21%)	2 (8%)	2 (9%)	6 (22%)
Development, validation, and incremental value study	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	2 (6%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)
Predesigned validation study	1 (6%)	1 (9%)	0 (0%)	0 (0%)	1 (3%)	1 (3%)	0 (0%)	5 (21%)	1 (4%)	10 (37%)
Validation and incremental value	4 (25%)	1 (9%)	0 (0%)	3 (13%)	2 (7%)	3 (9%)	1 (7%)	1 (4%)	9 (39%)	3 (11%)
Study design										
Existing cohort	12 (75%)	4 (36%)	3 (14%)	20 (87%)	20 (67%)	1 (3%)	2 (14%)	4 (17%)	7 (30%)	1 (4%)
Prospective cohort	1 (6%)	2 (18%)	0 (0%)	0 (0%)	1 (3%)	3 (9%)	0 (0%)	6 (25%)	4 (17%)	17 (63%)
Existing RCT	0 (0%)	1 (9%)	0 (0%)	0 (0%)	3 (10%)	0 (0%)	0 (0%)	0 (0%)	2 (9%)	0 (0%)
Existing registry	3 (19%)	4 (36%)	19 (86%)	1 (4%)	2 (7%)	30 (88%)	12 (86%)	14 (58%)	10 (43%)	9 (33%)
Case-control	0 (0%)	0 (0%)	0 (0%)	0 (0%)	4 (13%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)
Other	0 (0%)	0 (0%)	0 (0%)	2 (9%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)
Validation by independent investigators										
No	6 (38%)	3 (27%)	3 (14%)	10 (43%)	17 (57%)	2 (6%)	0 (0%)	7 (29%)	1 (4%)	2 (7%)
Yes	10 (62%)	8 (73%)	19 (86%)	13 (57%)	13 (43%)	32 (94%)	14 (100%)	17 (71%)	22 (96%)	25 (93%)
Comparability of eligibility criteria										
Narrower	6 (38%)	2 (18%)	18 (82%)	18 (78%)	28 (93%)	22 (65%)	0 (0%)	4 (17%)	4 (17%)	20 (74%)
Comparable	4 (25%)	0 (0%)	2 (9%)	3 (13%)	2 (7%)	5 (15%)	0 (0%)	1 (4%)	3 (13%)	7 (26%)
Mixture	5 (31%)	9 (82%)	0 (0%)	2 (9%)	0 (0%)	3 (9%)	0 (0%)	16 (67%)	11 (48%)	0 (0%)
Broader	1 (6%)	0 (0%)	2 (9%)	0 (0%)	0 (0%)	2 (6%)	0 (0%)	3 (12%)	5 (22%)	0 (0%)
Non-overlapping	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	14 (100%)	0 (0%)	0 (0%)	0 (0%)
Unclear	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	2 (6%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)
Comparability of age eligibility criteria										
Narrower	1 (6%)	2 (18%)	0 (0%)	17 (74%)	9 (30%)	13 (38%)	10 (71%)	2 (8%)	1 (4%)	17 (63%)
Comparable	15 (94%)	0 (0%)	22 (100%)	0 (0%)	1 (3%)	21 (62%)	4 (29%)	15 (62%)	4 (17%)	4 (15%)
Mixture	0 (0%)	0 (0%)	0 (0%)	6 (26%)	16 (53%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)
Broader	0 (0%)	9 (82%)	0 (0%)	0 (0%)	4 (13%)	0 (0%)	0 (0%)	7 (29%)	18 (78%)	6 (22%)
Setting										
Primary care	3 (19%)	0 (0%)	0 (0%)	7 (30%)	3 (10%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)
Secondary care	12 (75%)	9 (82%)	16 (84%)	0 (0%)	5 (17%)	18 (82%)	6 (75%)	18 (90%)	12 (86%)	4 (40%)
Tertiary care	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)

	ABCD2	ESRS	EuroSCORE	Framingham	FRAX	ISS	MELD	PSI	RCRI	SAPS 3
Population based	0 (0%)	1 (9%)	0 (0%)	15 (65%)	17 (59%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)
Screening	0 (0%)	0 (0%)	0 (0%)	1 (4%)	1 (3%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)
Mixed	1 (6%)	1 (9%)	2 (11%)	0 (0%)	0 (0%)	4 (18%)	1 (12%)	2 (10%)	2 (14%)	2 (20%)
Unclear	0 (0%)	0 (0%)	1 (5%)	0 (0%)	3 (10%)	0 (0%)	1 (12%)	0 (0%)	0 (0%)	4 (40%)
Comparability of setting										
Comparable	1 (6%)	0 (0%)	16 (73%)	15 (65%)	17 (57%)	18 (53%)	6 (43%)	18 (75%)	9 (39%)	4 (15%)
Narrower	15 (94%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)
Broader	0 (0%)	0 (0%)	2 (9%)	0 (0%)	0 (0%)	4 (12%)	1 (7%)	2 (8%)	2 (9%)	2 (7%)
Non-overlapping	0 (0%)	11 (100%)	3 (14%)	8 (35%)	10 (33%)	12 (35%)	6 (43%)	4 (17%)	12 (52%)	17 (63%)
Unclear	0 (0%)	0 (0%)	1 (5%)	0 (0%)	3 (10%)	0 (0%)	1 (7%)	0 (0%)	0 (0%)	4 (15%)
Continent										
Africa	0 (0%)	0 (0%)	1 (5%)	0 (0%)	0 (0%)	1 (3%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)
Asia	0 (0%)	2 (18%)	4 (18%)	2 (9%)	4 (13%)	7 (21%)	1 (7%)	3 (12%)	2 (9%)	5 (19%)
Australia	0 (0%)	0 (0%)	1 (5%)	0 (0%)	5 (17%)	2 (6%)	0 (0%)	3 (12%)	1 (4%)	1 (4%)
Europe	8 (50%)	7 (64%)	10 (45%)	10 (43%)	9 (30%)	7 (21%)	7 (50%)	11 (46%)	13 (57%)	9 (33%)
North America	8 (50%)	1 (9%)	5 (23%)	11 (48%)	11 (37%)	17 (50%)	3 (21%)	6 (25%)	6 (26%)	3 (11%)
South America	0 (0%)	0 (0%)	1 (5%)	0 (0%)	0 (0%)	0 (0%)	3 (21%)	0 (0%)	0 (0%)	9 (33%)
Combination	0 (0%)	1 (9%)	0 (0%)	0 (0%)	1 (3%)	0 (0%)	0 (0%)	1 (4%)	1 (4%)	0 (0%)
Comparability of continent										
Comparable	0 (0%)	0 (0%)	10 (45%)	11 (48%)	0 (0%)	17 (50%)	3 (21%)	6 (25%)	6 (26%)	0 (0%)
Narrower	16 (100%)	8 (73%)	0 (0%)	0 (0%)	30 (100%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	27 (100%)
Broader	0 (0%)	1 (9%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	1 (4%)	1 (4%)	0 (0%)
Non-overlapping	0 (0%)	2 (18%)	12 (55%)	12 (52%)	0 (0%)	17 (50%)	11 (79%)	17 (71%)	16 (70%)	0 (0%)
Number of centres										
Single	9 (56%)	4 (36%)	12 (55%)	3 (13%)	5 (17%)	18 (53%)	12 (86%)	14 (58%)	17 (74%)	13 (48%)
Multiple	6 (38%)	7 (64%)	9 (41%)	6 (26%)	9 (30%)	15 (44%)	2 (14%)	10 (42%)	6 (26%)	14 (52%)
Population based	1 (6%)	0 (0%)	0 (0%)	12 (52%)	15 (50%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)
Unclear	0 (0%)	0 (0%)	1 (5%)	2 (9%)	1 (3%)	1 (3%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)
Changes made to predictors										
No	16 (100%)	10 (91%)	13 (59%)	23 (100%)	20 (67%)	21 (62%)	12 (86%)	13 (54%)	17 (74%)	26 (96%)
Yes	0 (0%)	0 (0%)	5 (23%)	0 (0%)	10 (33%)	0 (0%)	0 (0%)	5 (21%)	2 (9%)	0 (0%)
Unclear	0 (0%)	1 (9%)	4 (18%)	0 (0%)	0 (0%)	13 (38%)	2 (14%)	6 (25%)	4 (17%)	1 (4%)
Comparability of outcome definition										

	ABCD2	ESRS	EuroSCORE	Framingham	FRAX	ISS	MELD	PSI	RCRI	SAPS 3
No	8 (50%)	3 (27%)	9 (41%)	13 (57%)	16 (53%)	7 (21%)	14 (100%)	5 (21%)	5 (22%)	24 (89%)
Yes	3 (19%)	8 (73%)	11 (50%)	4 (17%)	13 (43%)	19 (56%)	0 (0%)	18 (75%)	18 (78%)	3 (11%)
Unclear	5 (31%)	0 (0%)	2 (9%)	6 (26%)	1 (3%)	8 (24%)	0 (0%)	1 (4%)	0 (0%)	0 (0%)
Outcome measurement method										
Self-reported	3 (19%)	8 (73%)	2 (9%)	15 (65%)	18 (60%)	0 (0%)	1 (7%)	6 (25%)	1 (4%)	2 (7%)
Clinician	6 (38%)	1 (9%)	2 (9%)	2 (9%)	1 (3%)	4 (12%)	0 (0%)	2 (8%)	10 (43%)	6 (22%)
Registry	3 (19%)	2 (18%)	5 (23%)	4 (17%)	8 (27%)	13 (38%)	3 (21%)	9 (38%)	6 (26%)	9 (33%)
Unclear	4 (25%)	0 (0%)	13 (59%)	2 (9%)	3 (10%)	17 (50%)	10 (71%)	7 (29%)	6 (26%)	10 (37%)
Similar outcome measurement for all patients										
Yes	9 (56%)	3 (27%)	11 (50%)	6 (26%)	14 (47%)	18 (53%)	2 (14%)	12 (50%)	13 (57%)	16 (59%)
No	1 (6%)	5 (45%)	1 (5%)	16 (70%)	3 (10%)	3 (9%)	0 (0%)	3 (12%)	4 (17%)	1 (4%)
Unclear	6 (38%)	3 (27%)	10 (45%)	1 (4%)	13 (43%)	13 (38%)	12 (86%)	9 (38%)	6 (26%)	10 (37%)
Method for handling of missing data										
Complete case analysis	4 (25%)	8 (73%)	3 (14%)	11 (48%)	19 (63%)	18 (53%)	7 (50%)	1 (4%)	5 (22%)	8 (30%)
Mean/median imputation	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	1 (4%)
Multiple imputation	1 (6%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)
NA	3 (19%)	1 (9%)	1 (5%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	1 (4%)	1 (4%)	2 (7%)
Other	0 (0%)	0 (0%)	1 (5%)	0 (0%)	2 (7%)	0 (0%)	1 (7%)	8 (33%)	0 (0%)	7 (26%)
Unclear	8 (50%)	2 (18%)	17 (77%)	12 (52%)	9 (30%)	16 (47%)	6 (43%)	14 (58%)	17 (74%)	9 (33%)
Handling of missing data										
Appropriate or <5% missing	8 (50%)	4 (36%)	2 (9%)	3 (13%)	4 (13%)	6 (18%)	4 (29%)	4 (17%)	7 (30%)	6 (22%)
Inappropriate	1 (6%)	5 (45%)	3 (14%)	8 (35%)	17 (57%)	13 (38%)	4 (29%)	8 (33%)	0 (0%)	12 (44%)
Unclear	7 (44%)	2 (18%)	17 (77%)	12 (52%)	9 (30%)	15 (44%)	6 (43%)	12 (50%)	16 (70%)	9 (33%)

Continuous variables

	ABCD2	ESRS	EuroSCORE	Framingham	FRAX	ISS	MELD	PSI	RCRI	SAPSIII
Year start recruitment	2002 (2000-2003) NR=0	2007 (2004-2007) NR=0	1998 (1995-2001) NR=1	1989 (1983-1994) NR=0	1994 (1990-1998) NR=3	1996 (1993-1998) NR=1	2000 (1998-2004) NR=0	2000 (1998-2002) NR=0	2000 (1994-2002) NR=4	2006 (2006-2007) NR=0
Year end recruitment	2005 (2003-2007) NR=0	2008 (2006-2008) NR=0	2002 (1999-2005) NR=1	1993 (1988-1998) NR=0	1997 (1993-2006) NR=8	2000 (1996-2003) NR=2	2006 (2004-2007) NR=0	2002 (2000-2003) NR=0	2002 (2000-2005) NR=4	2007 (2006-2009) NR=0
Percentage missings	0.95 (0.00-5.00) NR=7	5.12 (1.99-17.80) NR=2	6.40 (1.50-11.83) NR=18	4.90 (2.70-9.80) NR=18	30.25 (2.75-33.80) NR=16	9.05 (2.40-14.65) NR=20	4.05 (2.73-10.93) NR=8	0.52 (0.07-9.26) NR=18	1.00 (0.09-1.91) NR=16	5.85 (0.52-18.93) NR=15
Number of participants	304 (204-691) NR=0	1257 (712-2594) NR=0	1730 (873-4518) NR=2	2399 (928-4609) NR=0	2210 (889-6586) NR=0	2590 (960-20713) NR=0	418 (118-483) NR=0	730 (326-970) NR=1	496 (180-1480) NR=0	864 (485-1856) NR=0
Number of events	9 (3-18) NR=0	92 (60-134) NR=0	36 (13-87) NR=2	92 (72-160) NR=1	250 (86-581) NR=0	256 (113-1660) NR=2	49 (22-112) NR=0	54 (28-111) NR=1	31 (14-76) NR=0	180 (124-311) NR=1
Age mean	67.4 (64.1-70.0) NR=5	68.3 (67.1-71.5) NR=3	63.9 (62.5-65.2) NR=2	54.6 (50.9-58.3) NR=2	66.8 (63.0-71.3) NR=1	38.1 (32.4-41.3) NR=10	51.8 (49.1-53.0) NR=0	66.2 (64.0-69.3) NR=2	67.8 (66.0-71.9) NR=2	62.2 (60.8-64.8) NR=1
Age sd	13.8 (13.0-14.9) NR=5	12.4 (12.0-13.0) NR=1	9.3 (9.0-10.6) NR=8	7.3 (4.1-9.4) NR=0	8.3 (5.9-9.8) NR=0	20.9 (18.1-24.8) NR=2	10.0 (9.6-12.0) NR=1	17.8 (17.0-20.1) NR=3	10.0 (8.8-12.5) NR=4	17.0 (15.4-19.0) NR=3
Gender percentage men	47 (45-53) NR=2	57 (55-59) NR=1	77 (71-79) NR=1	100 (100-100) NR=0	0 (0-0) NR=0	71 (64-75) NR=11	68 (63-69) NR=0	57 (53-64) NR=1	67 (52-76) NR=0	59 (55-64) NR=0

Values represent median (IQR), number of missing values.

Figure S1: Associations between categorical variables and c-statistic

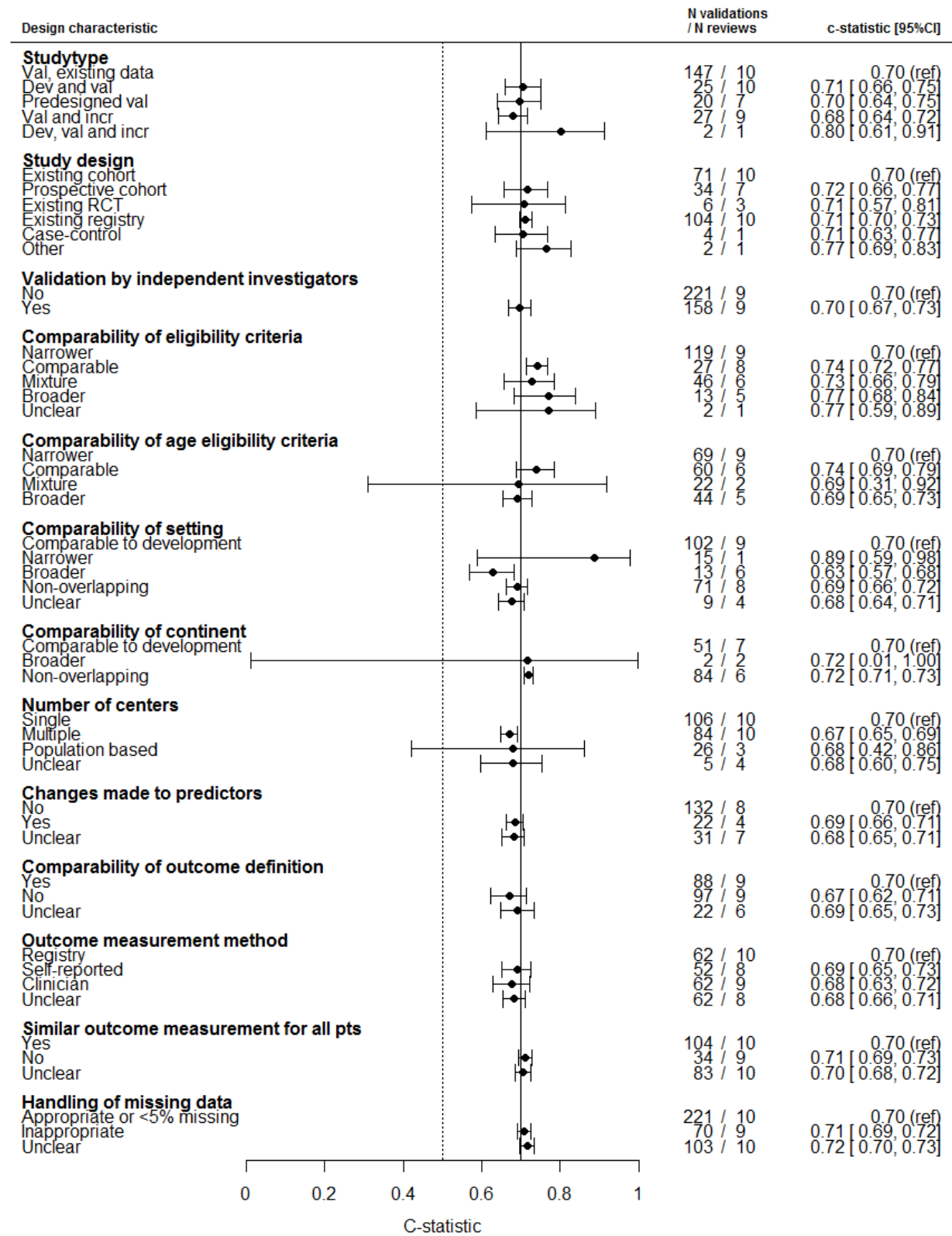
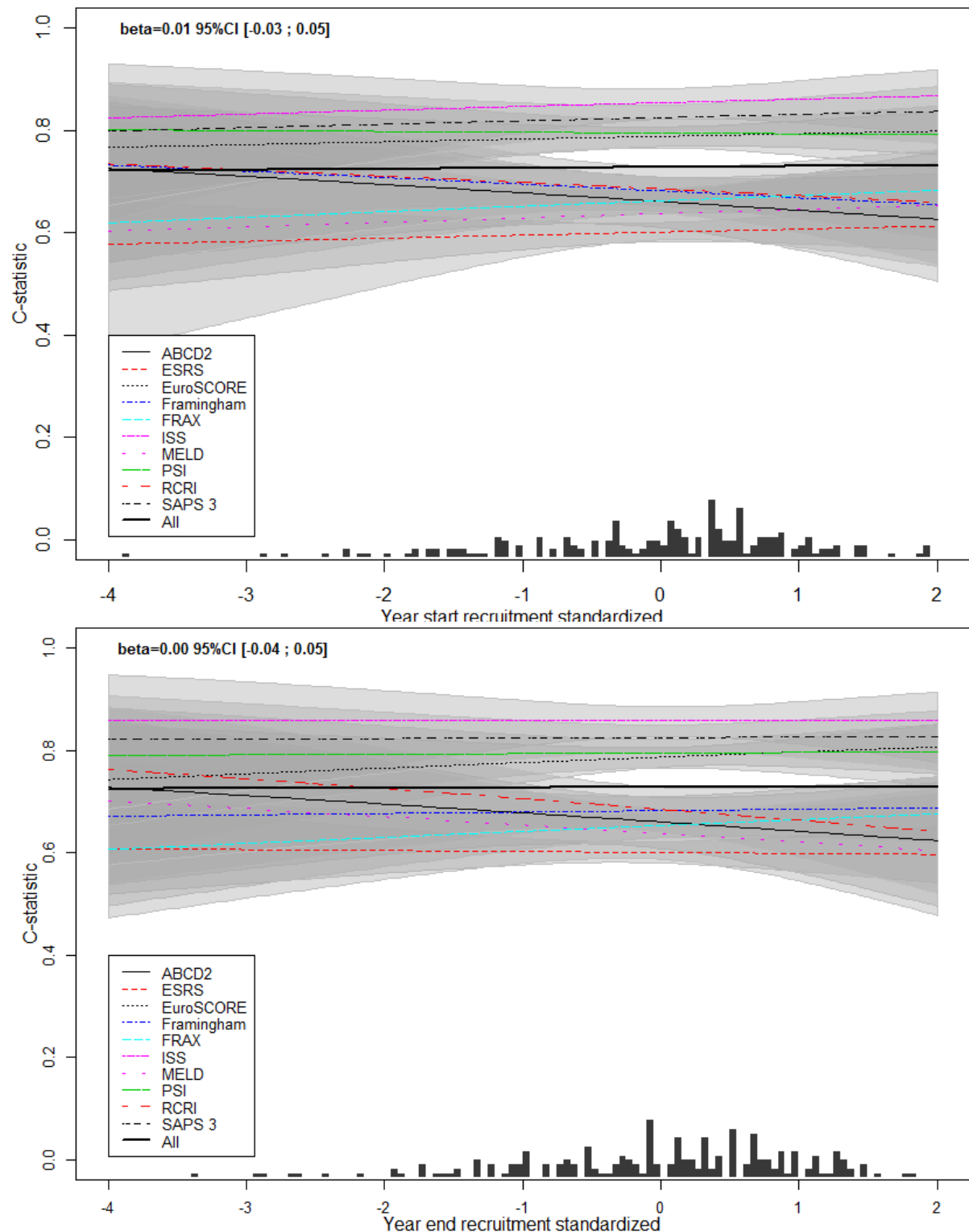
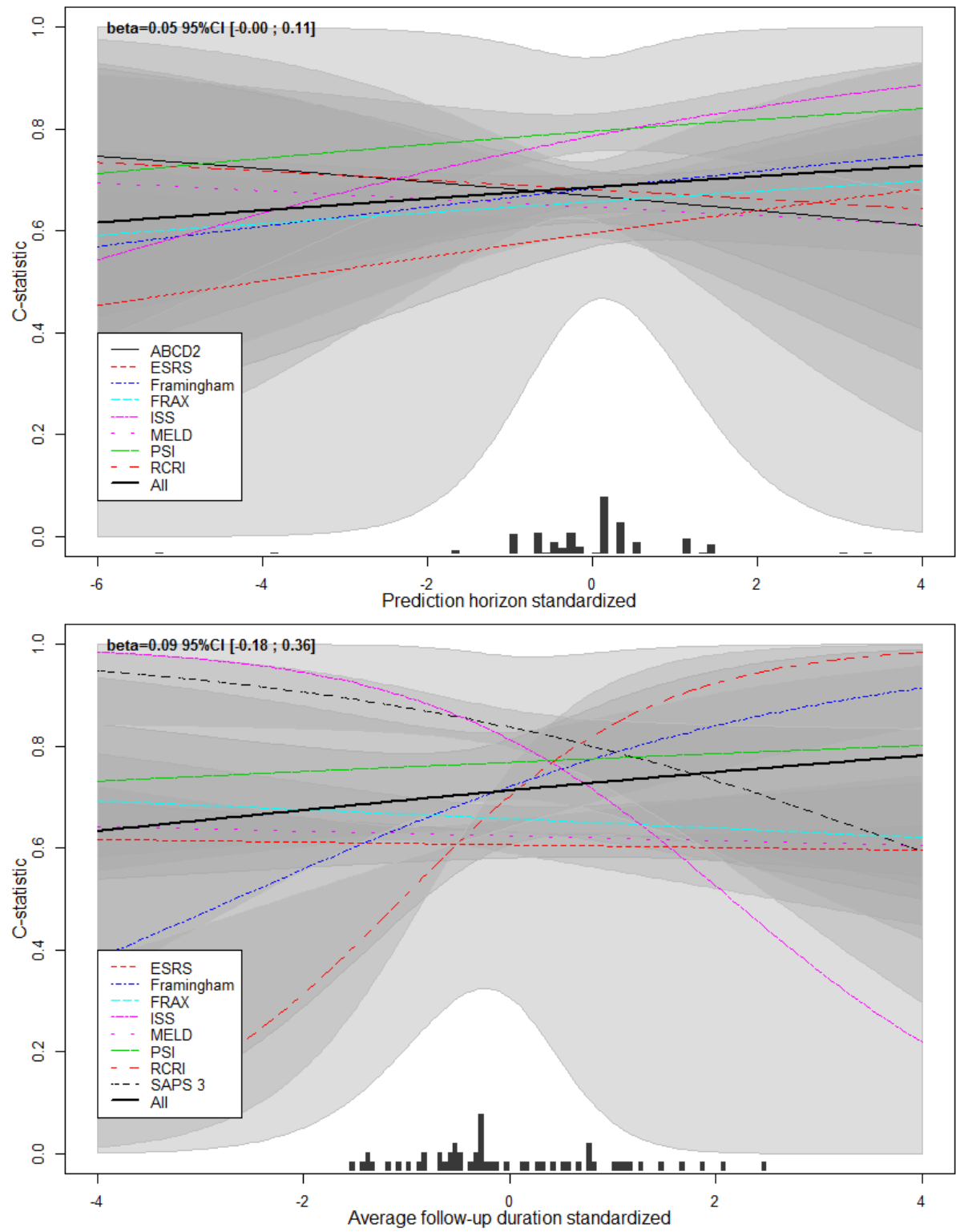
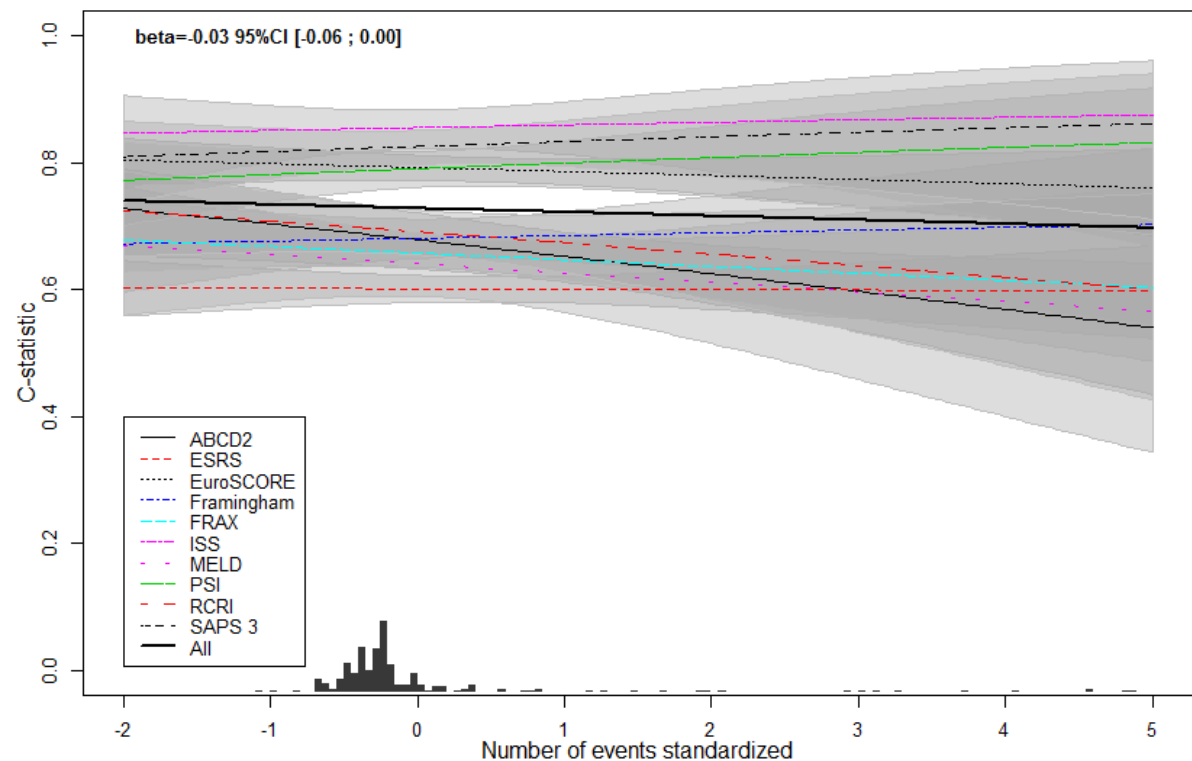
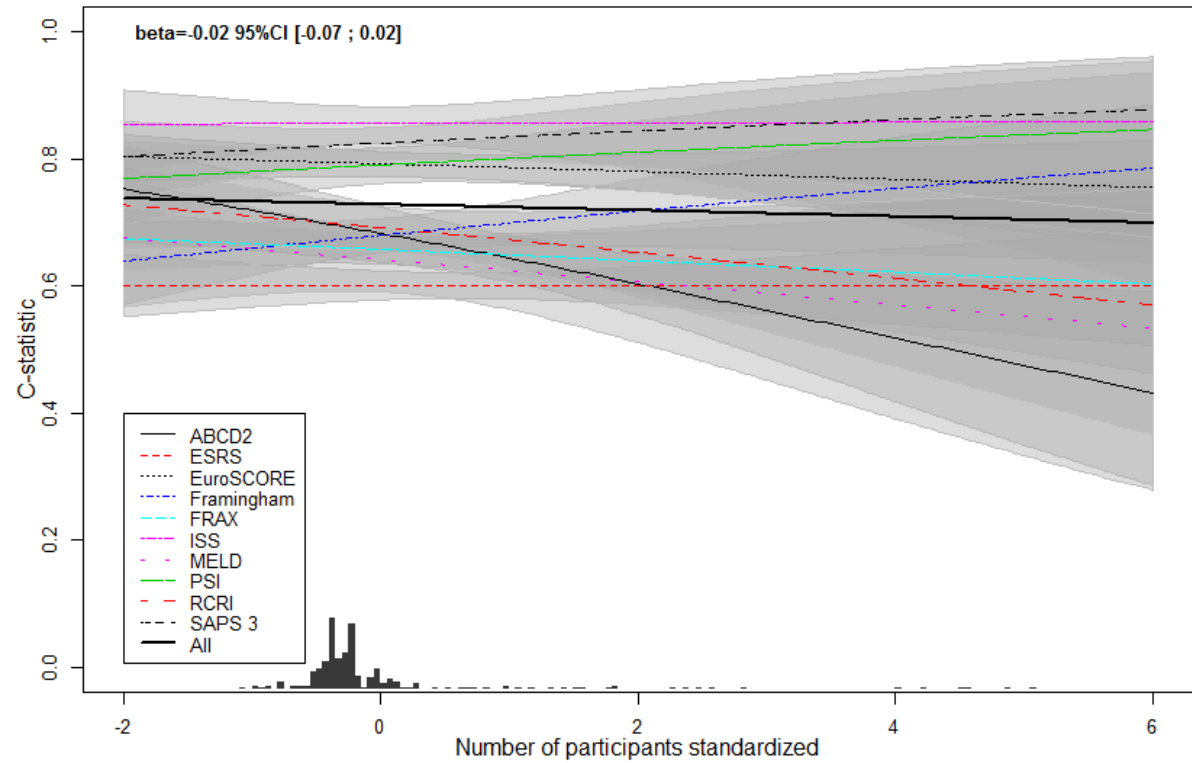
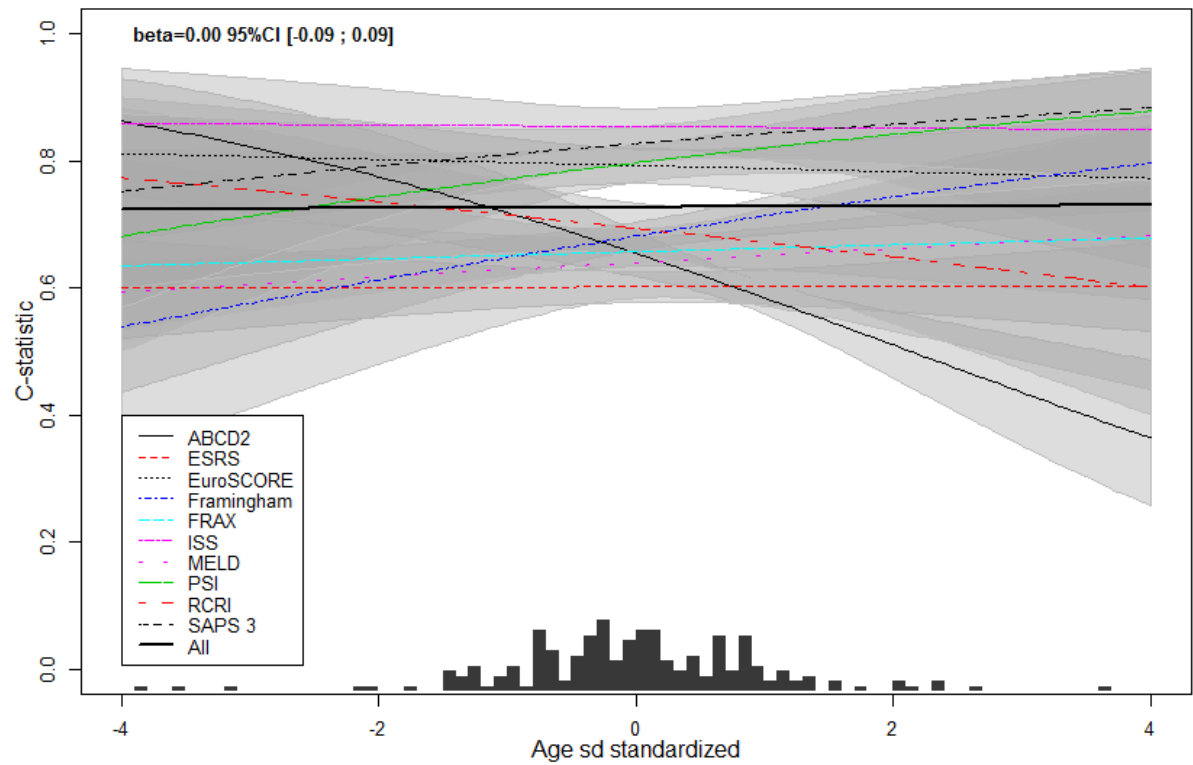
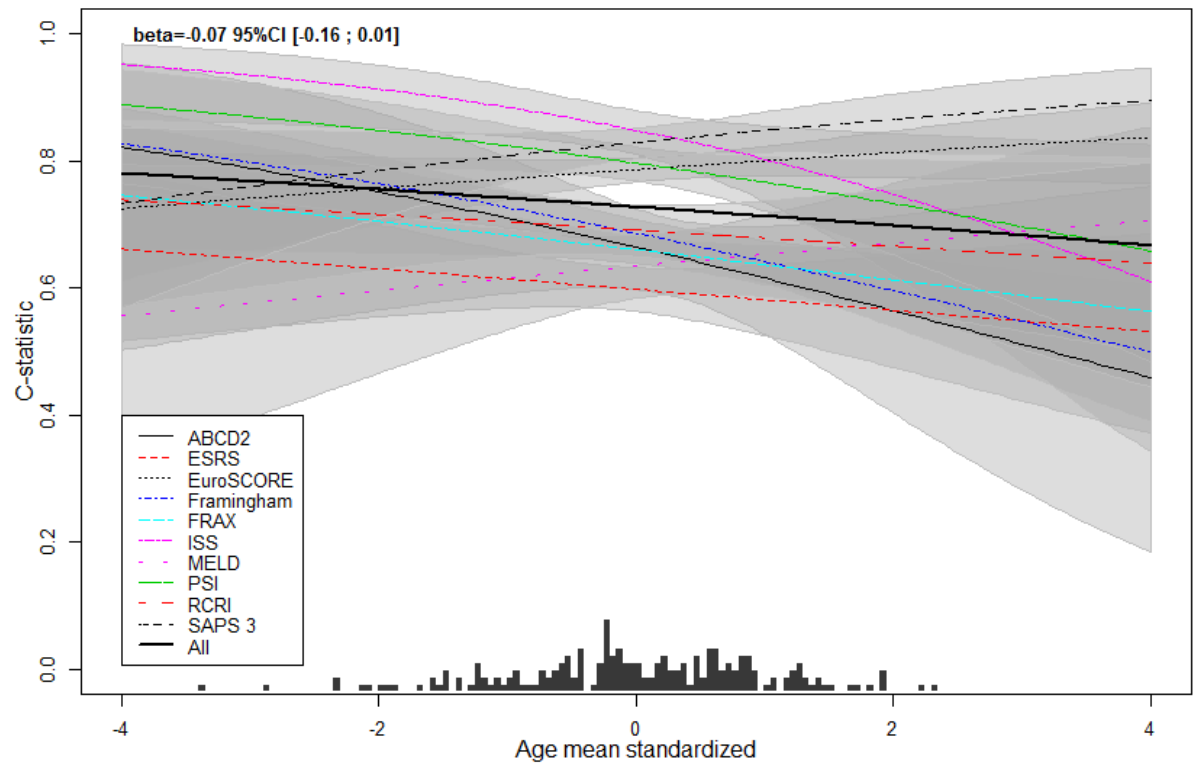


Figure S2: Associations between continuous variables and c-statistic









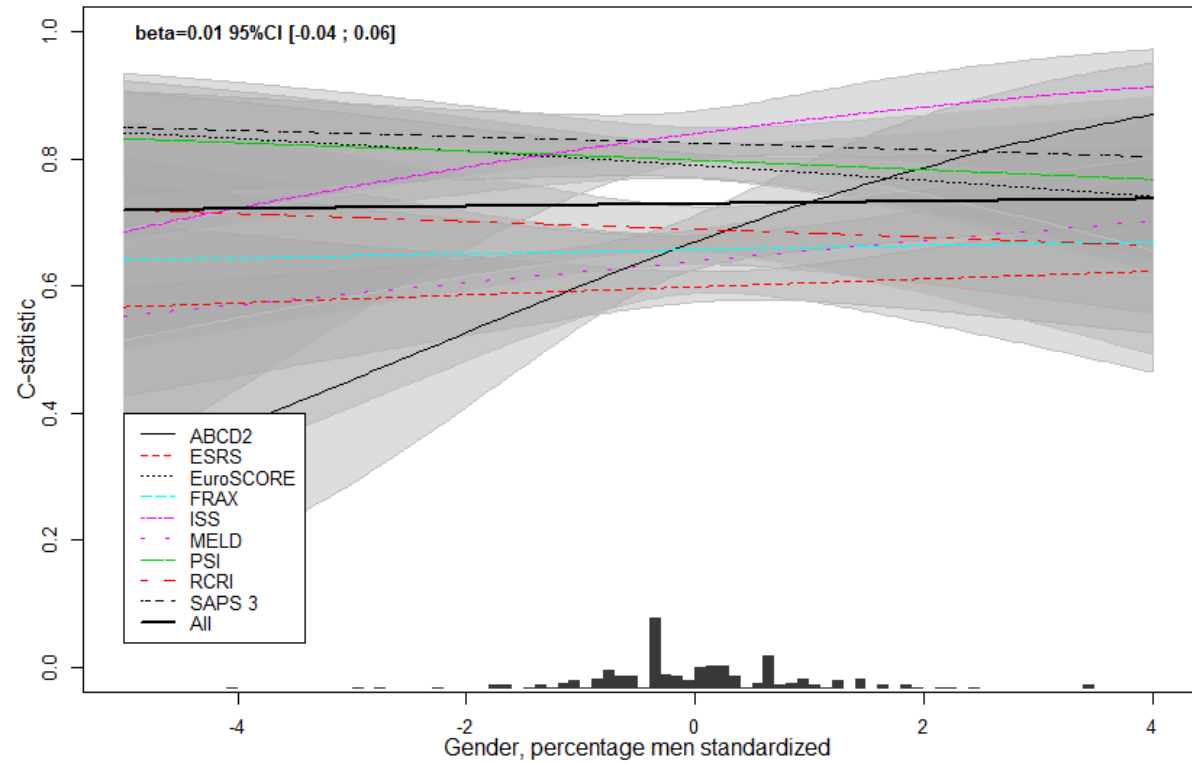
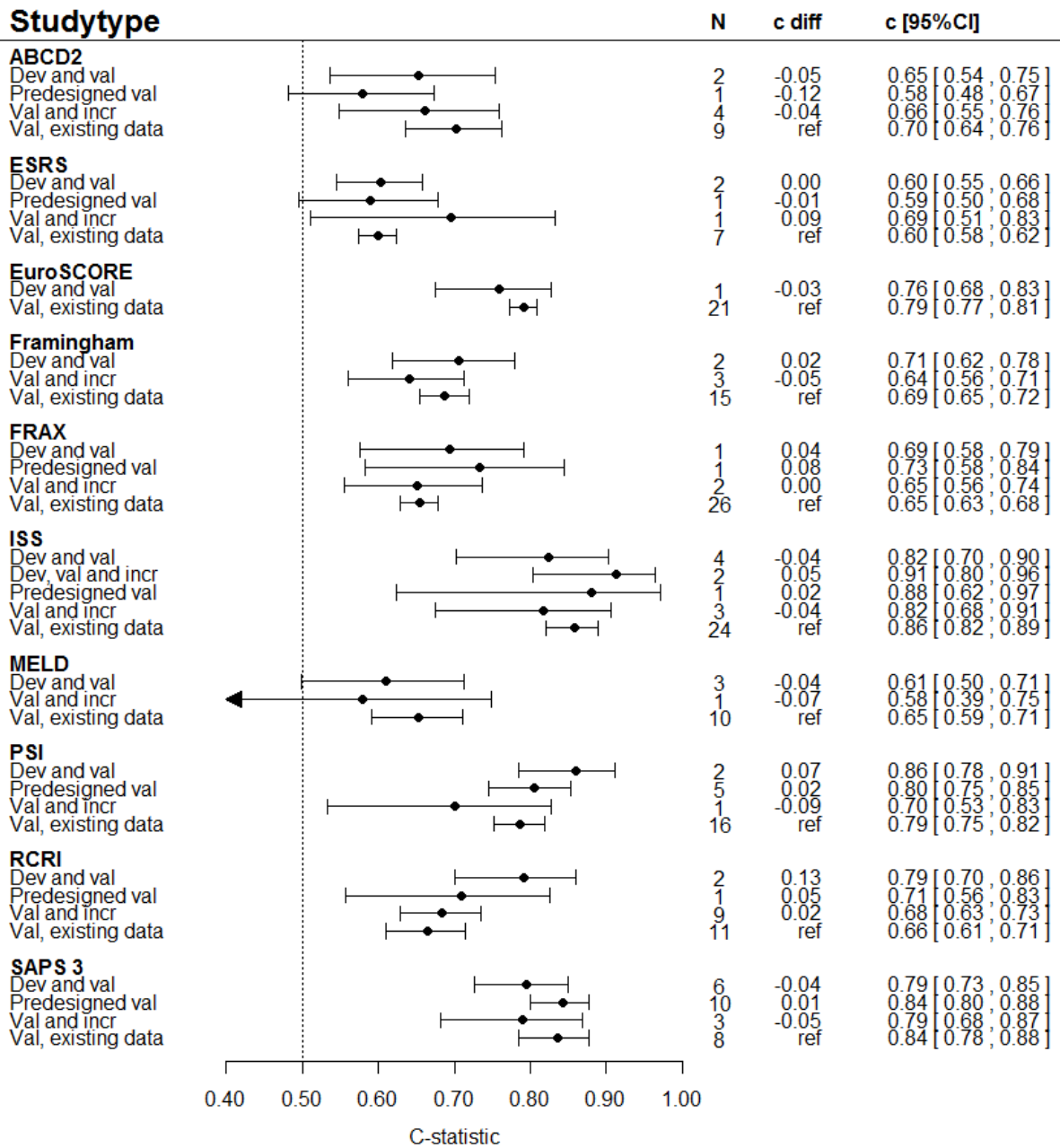
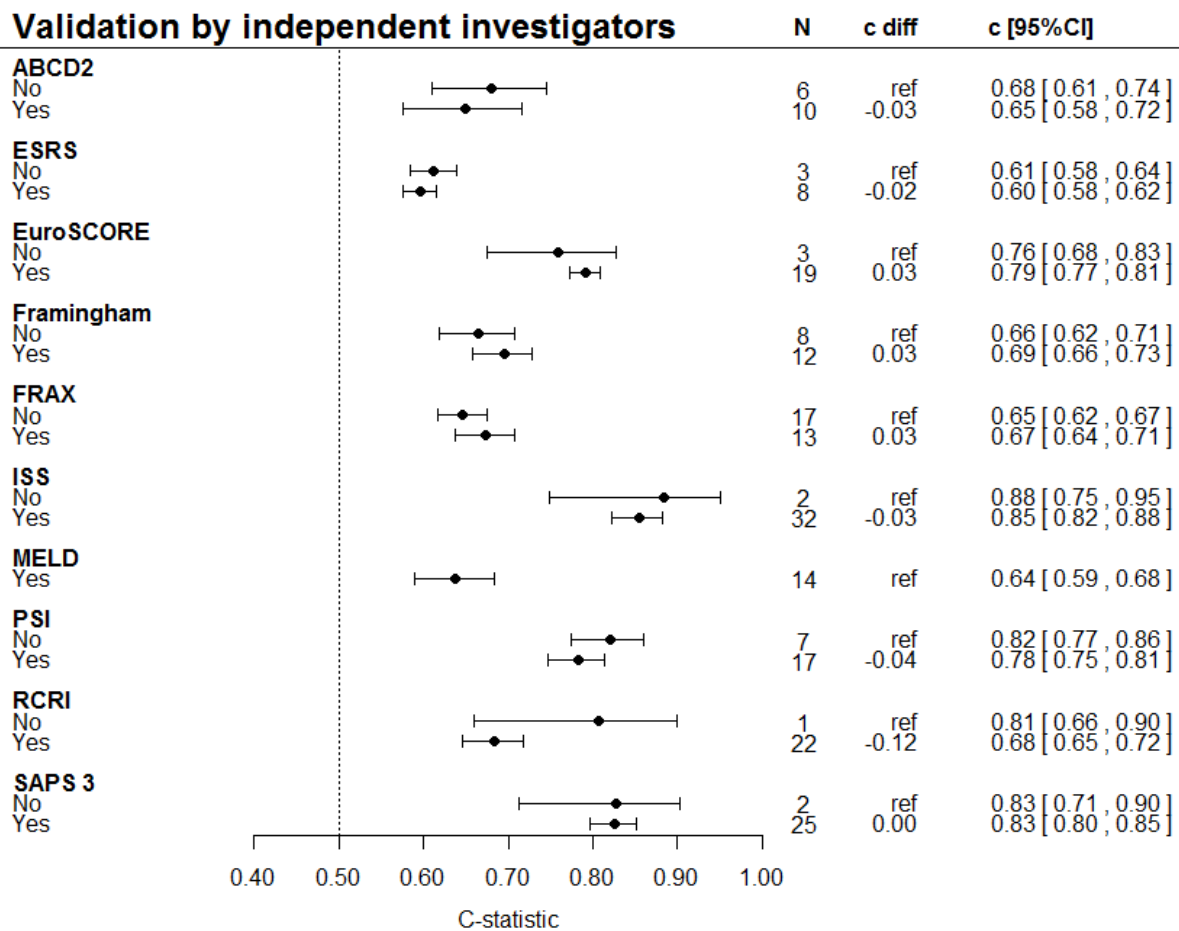
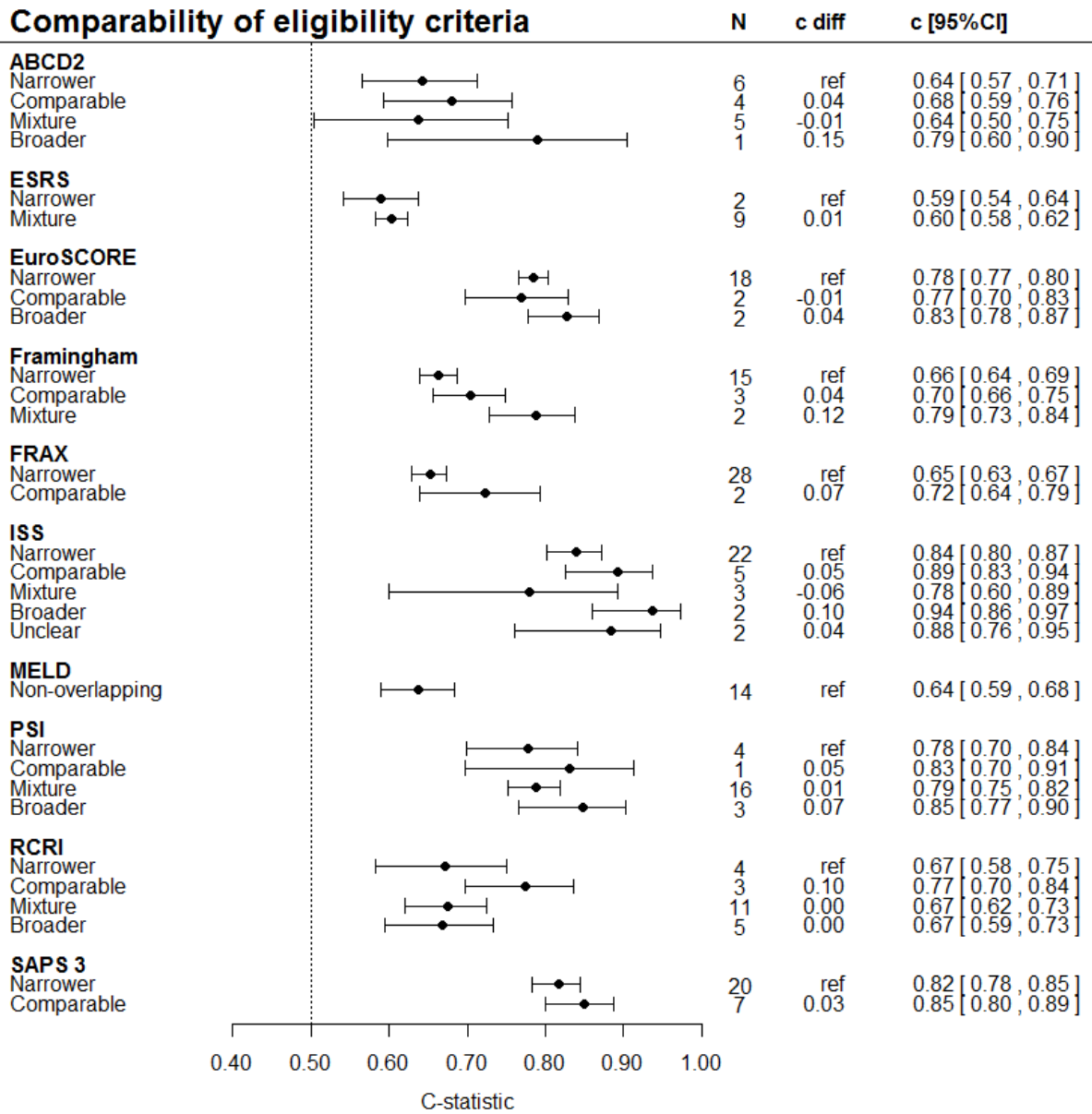
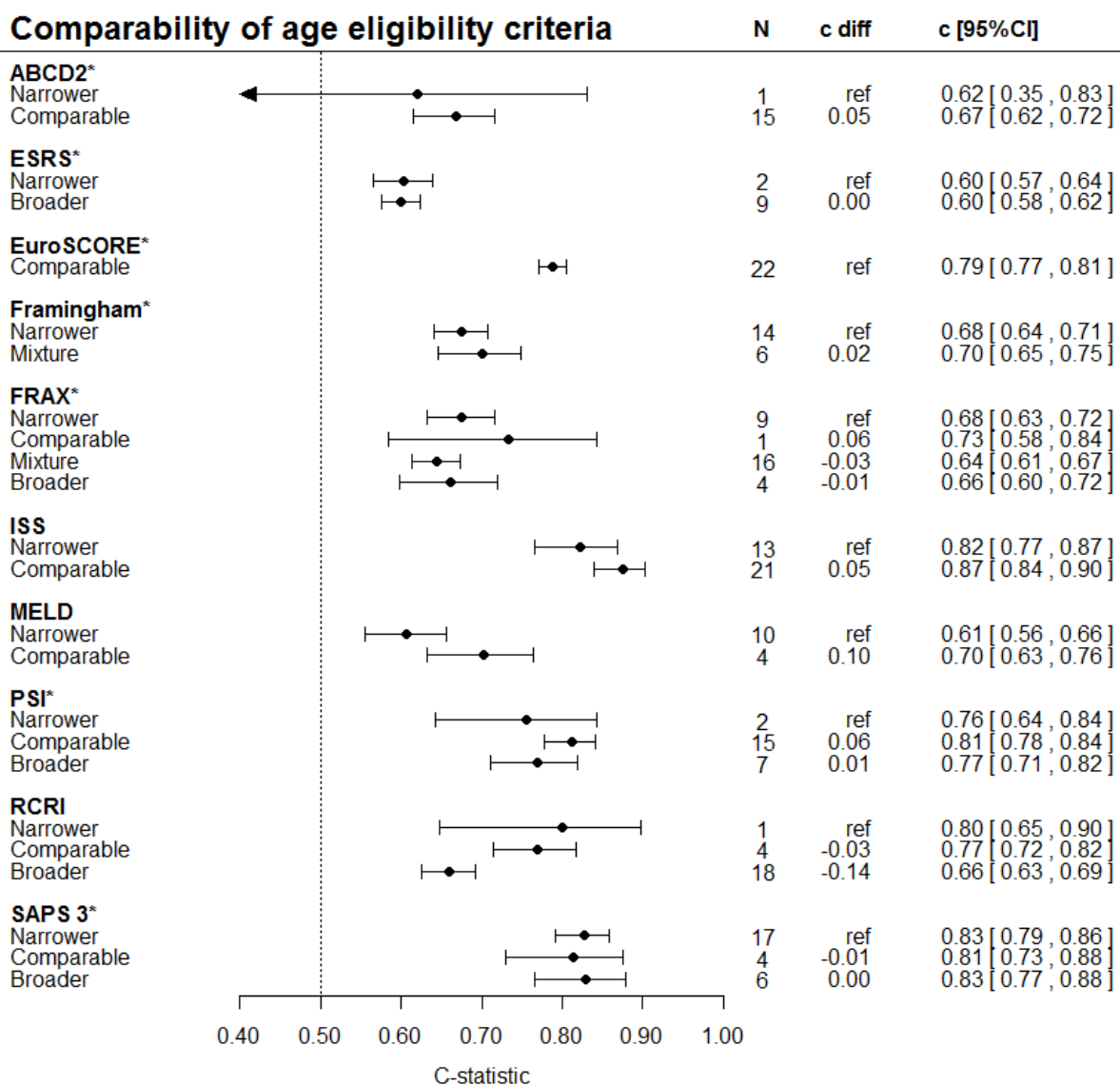


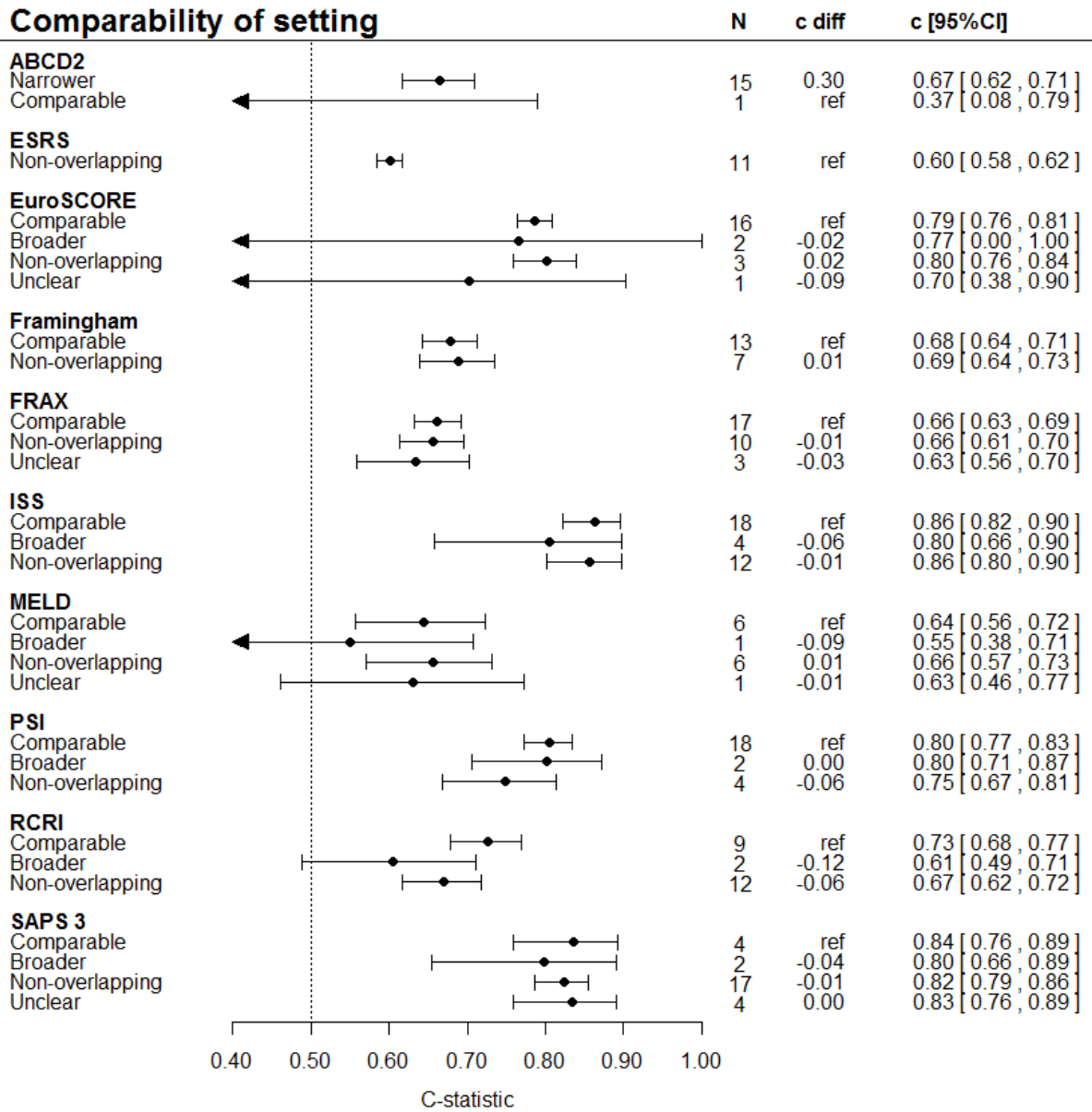
Figure S3: C-statistic in categories of study characteristics within each systematic review



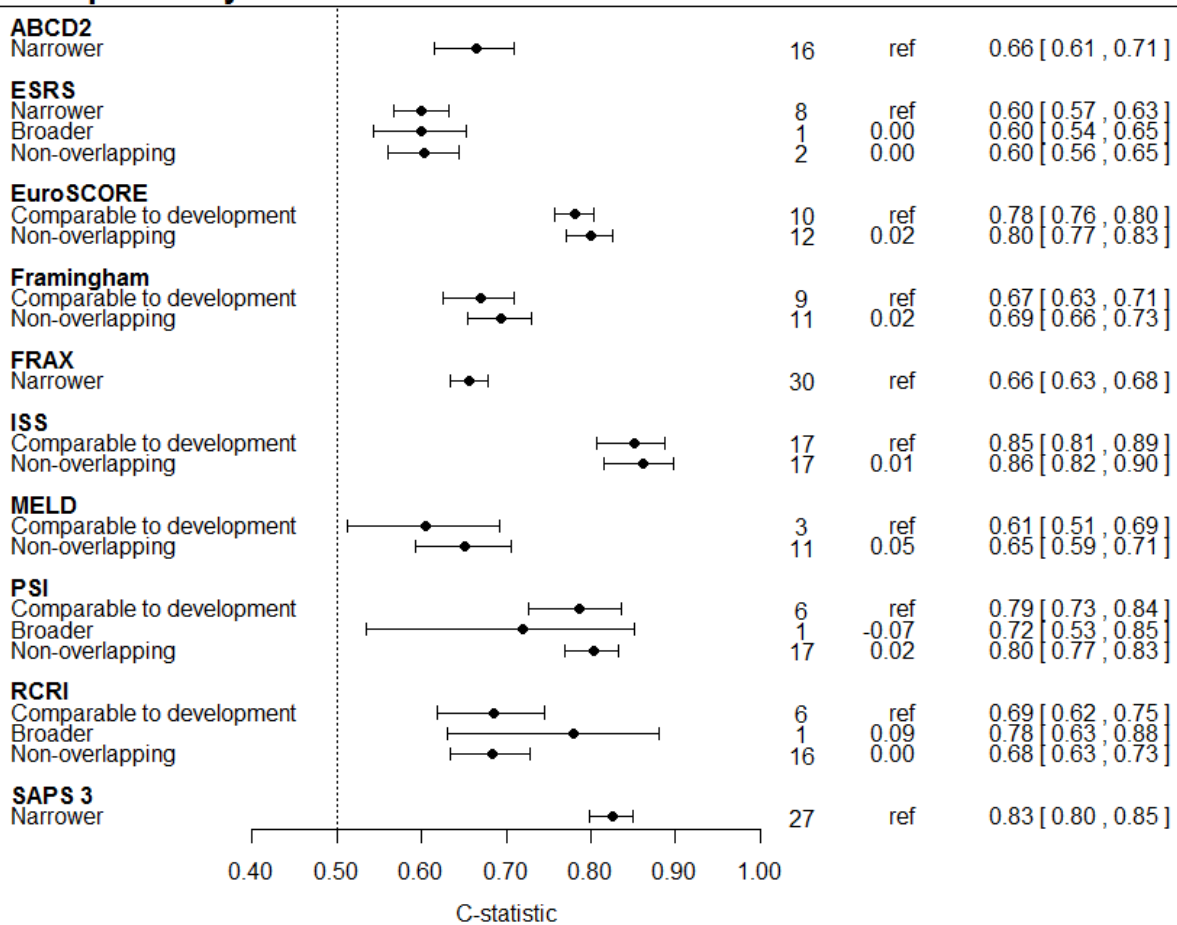


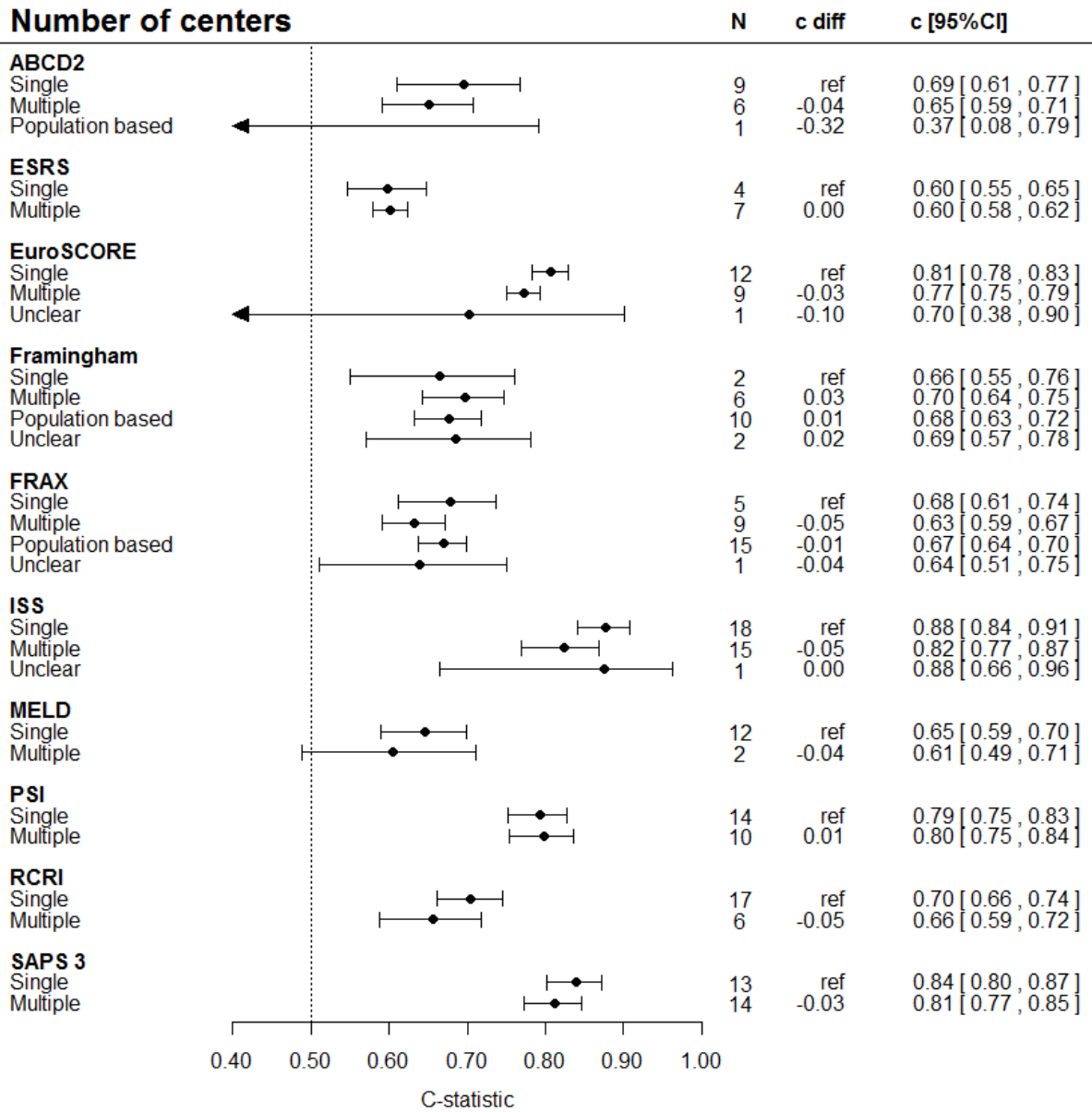


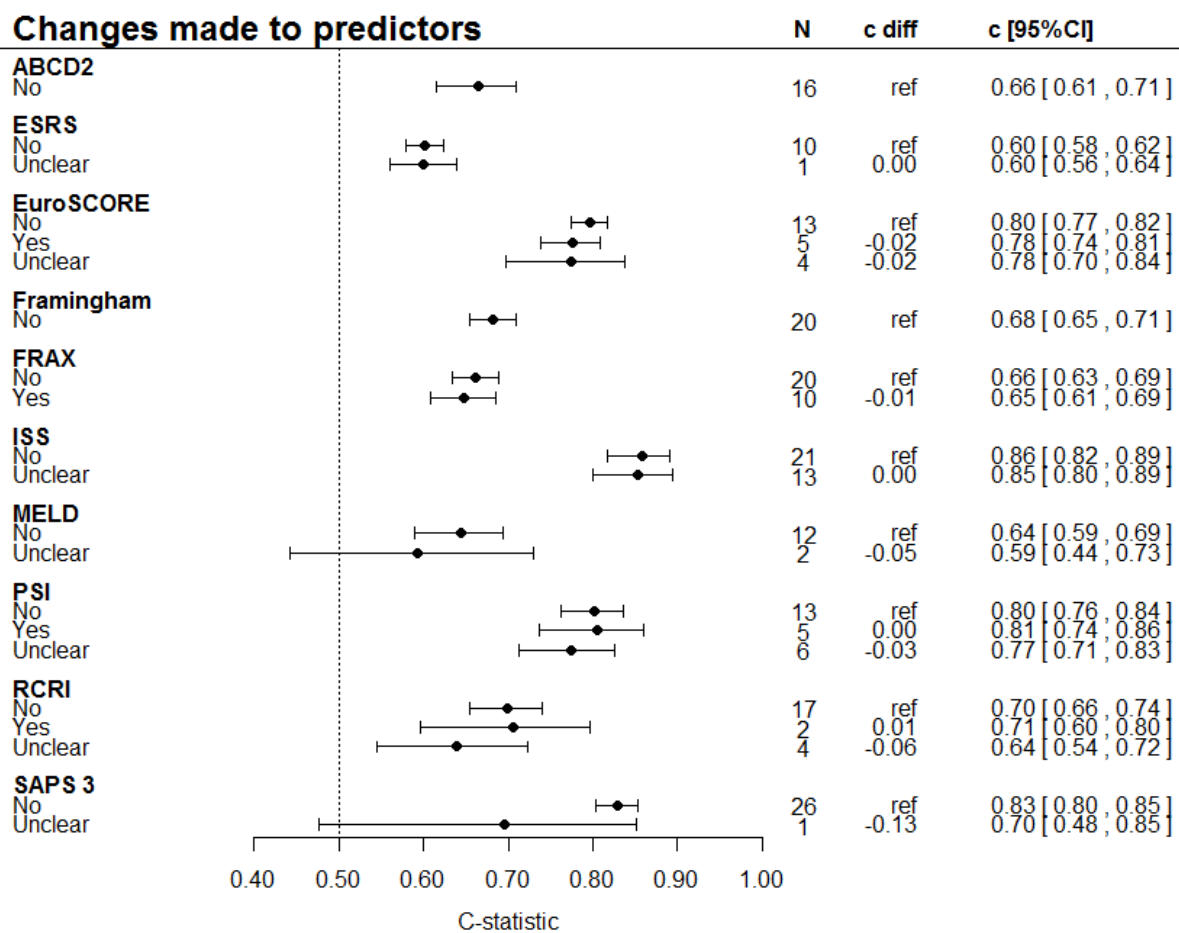


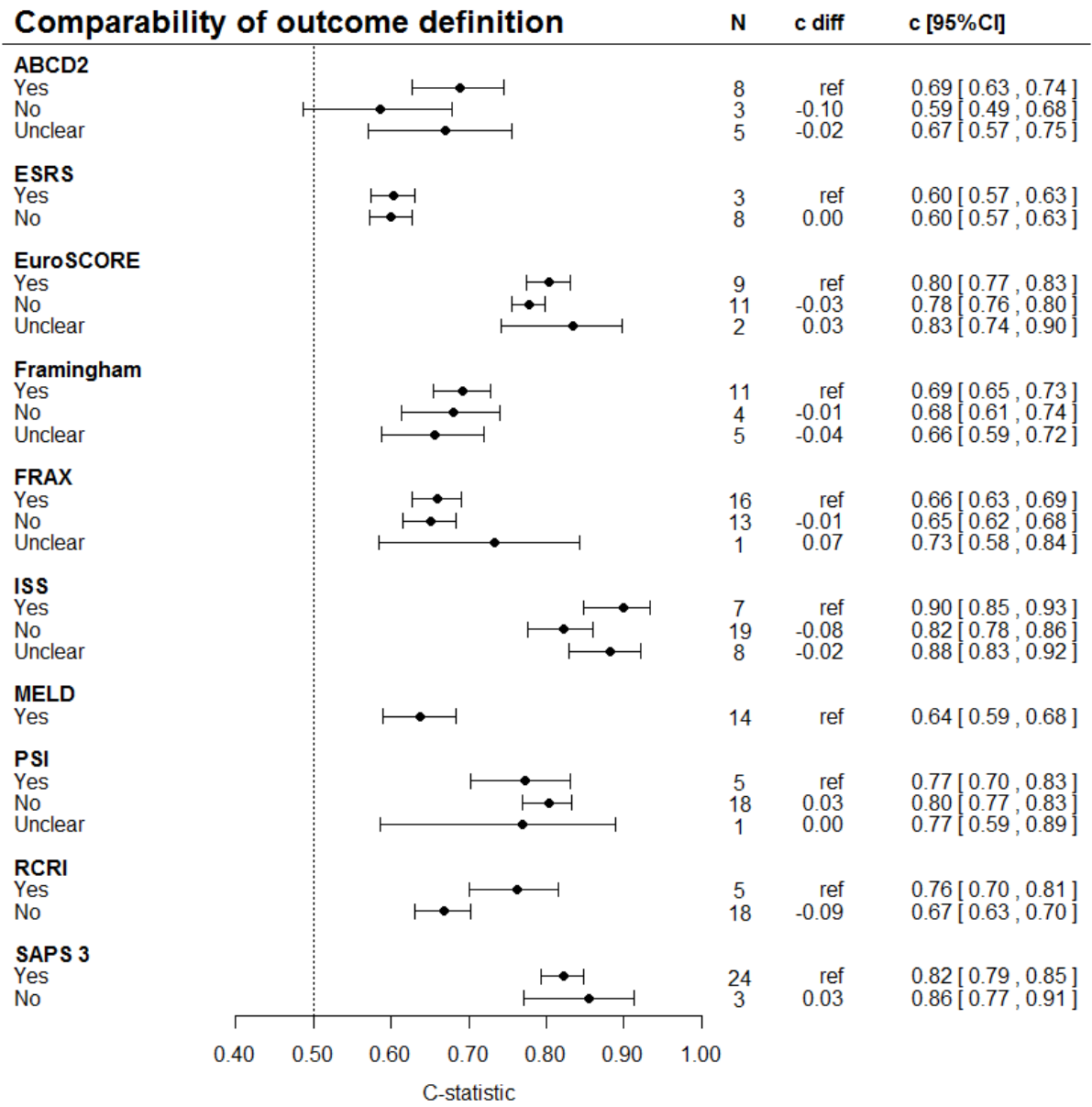


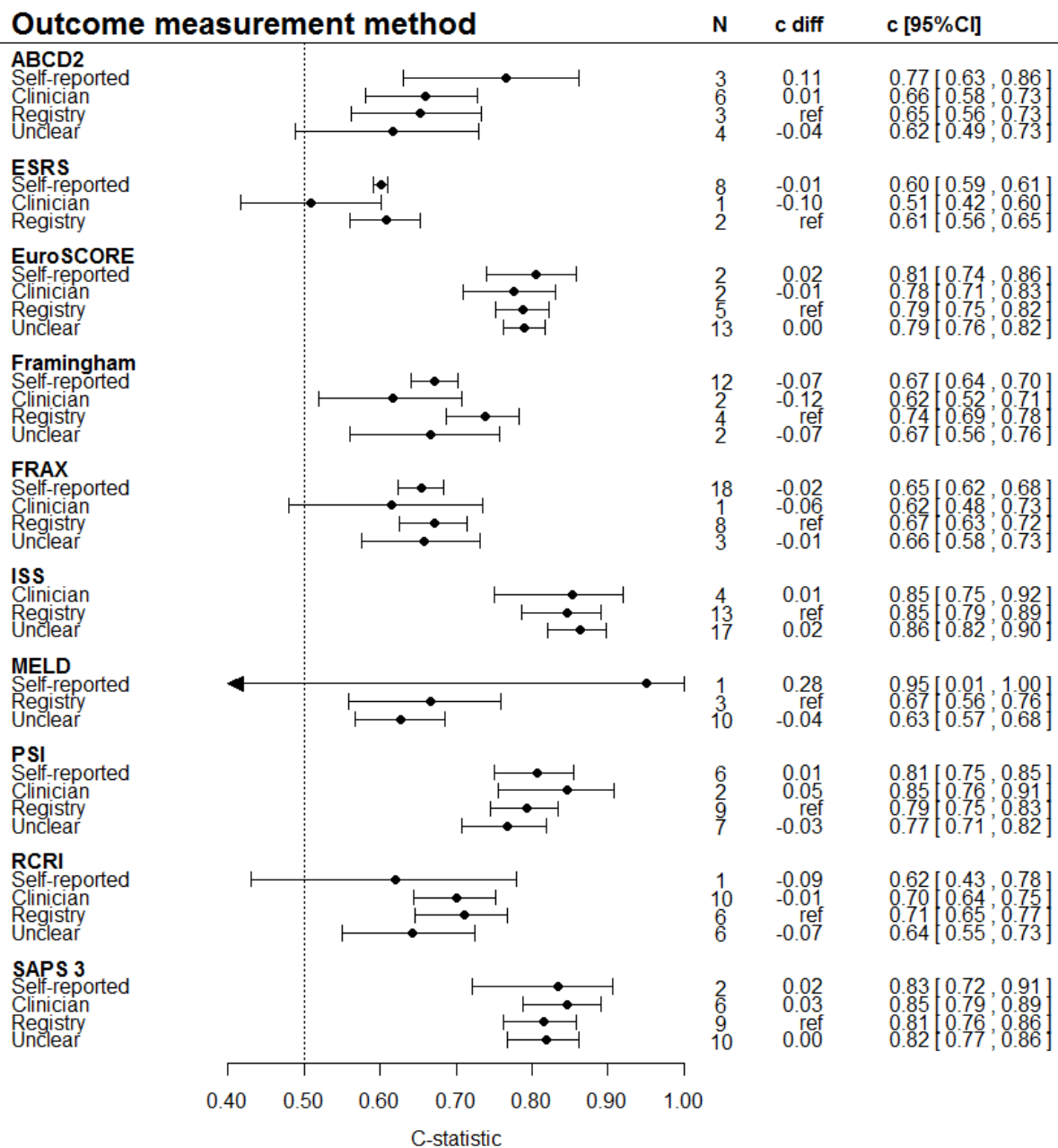
Comparability of continent

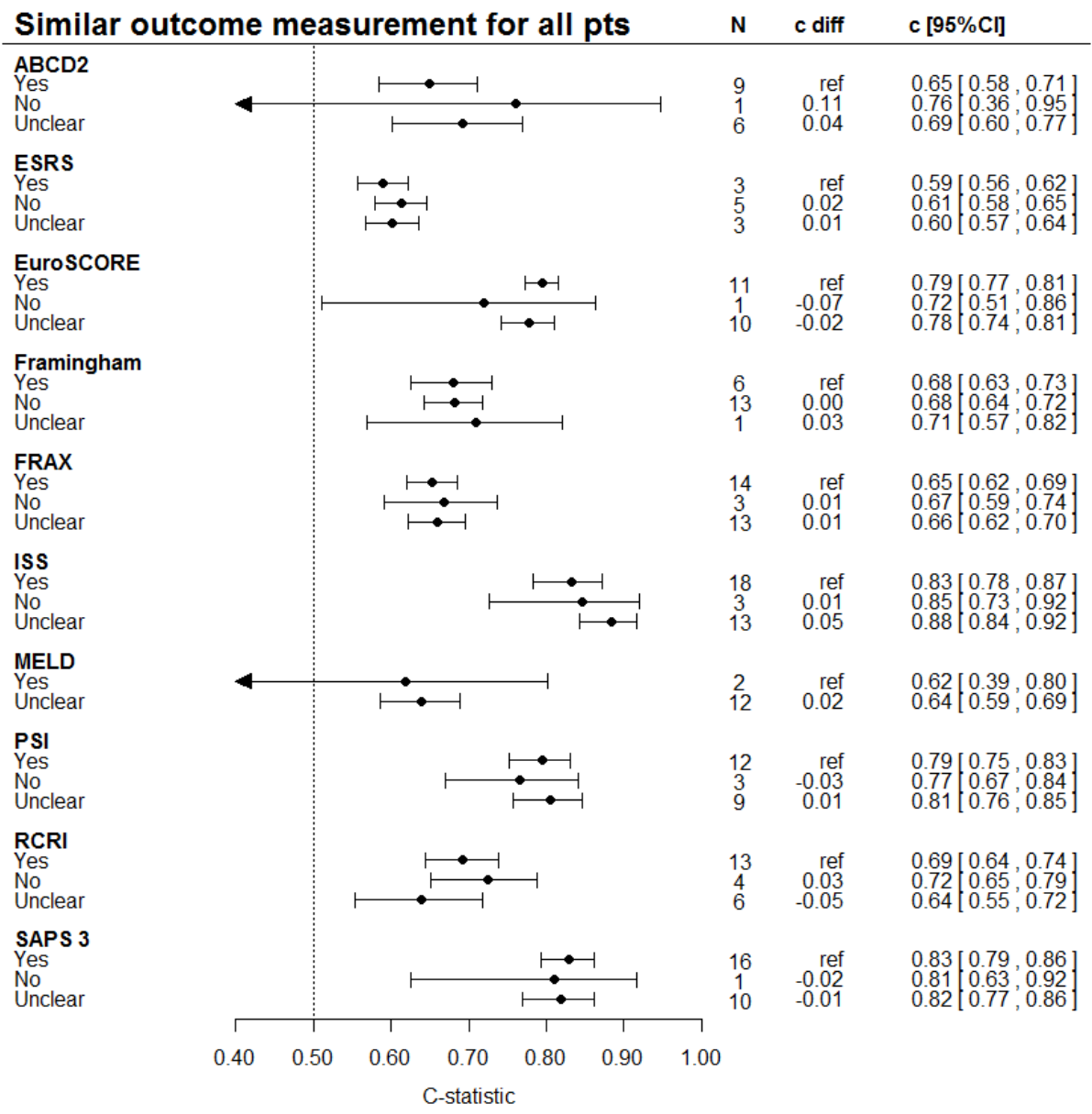


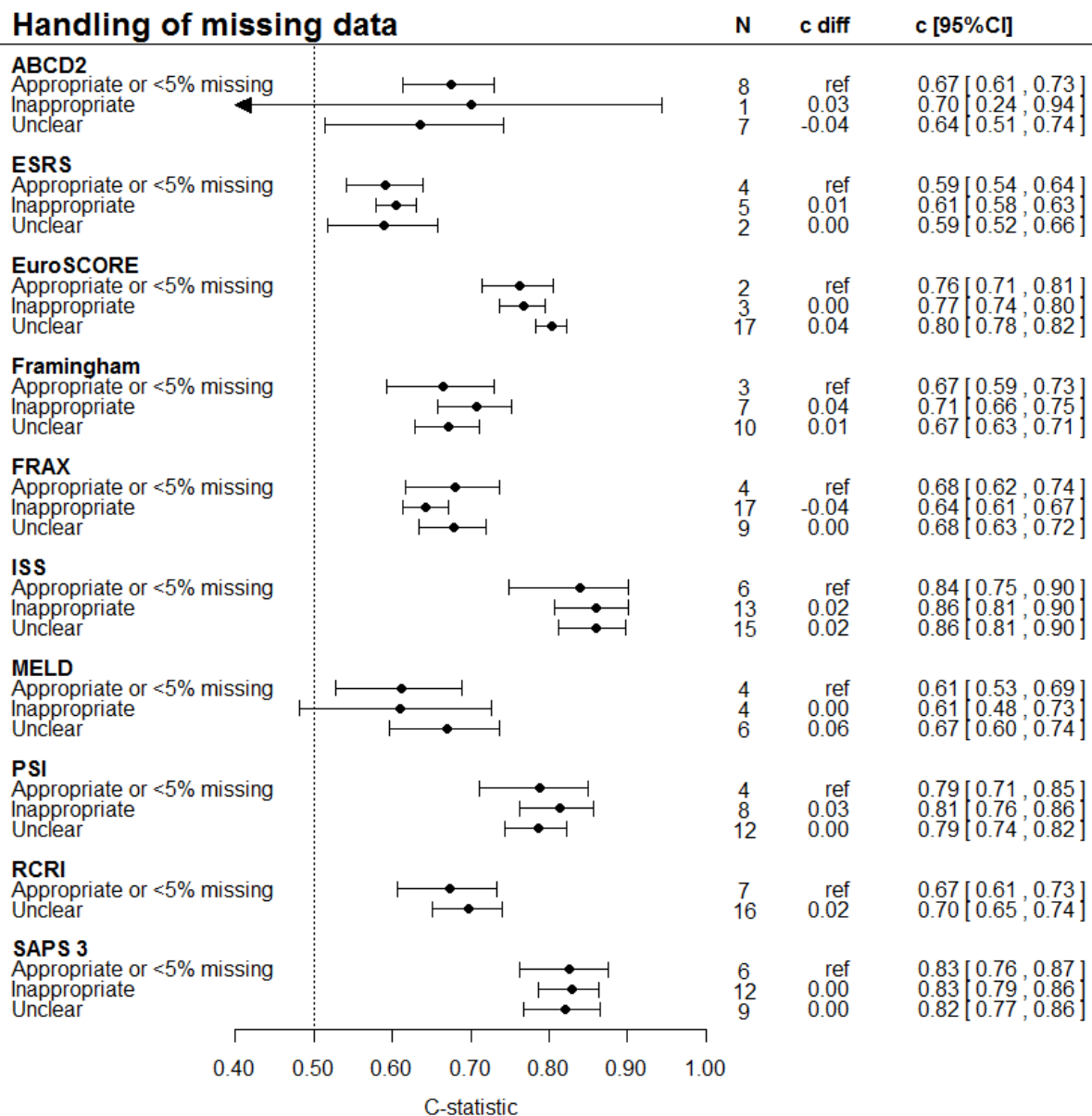












C-statistic for categories of study characteristics, pooled using univariable meta-regression analyses per systematic review. N represents the number of external validation studies in a specific category. C diff represents the difference in c-statistic with regard to a reference category (indicated with 'ref'). Dev: development, val: validation, incr: incremental value, pts: patients.

*Models contain age as predictor

Figure S4: Associations between categorical variables and OE ratio

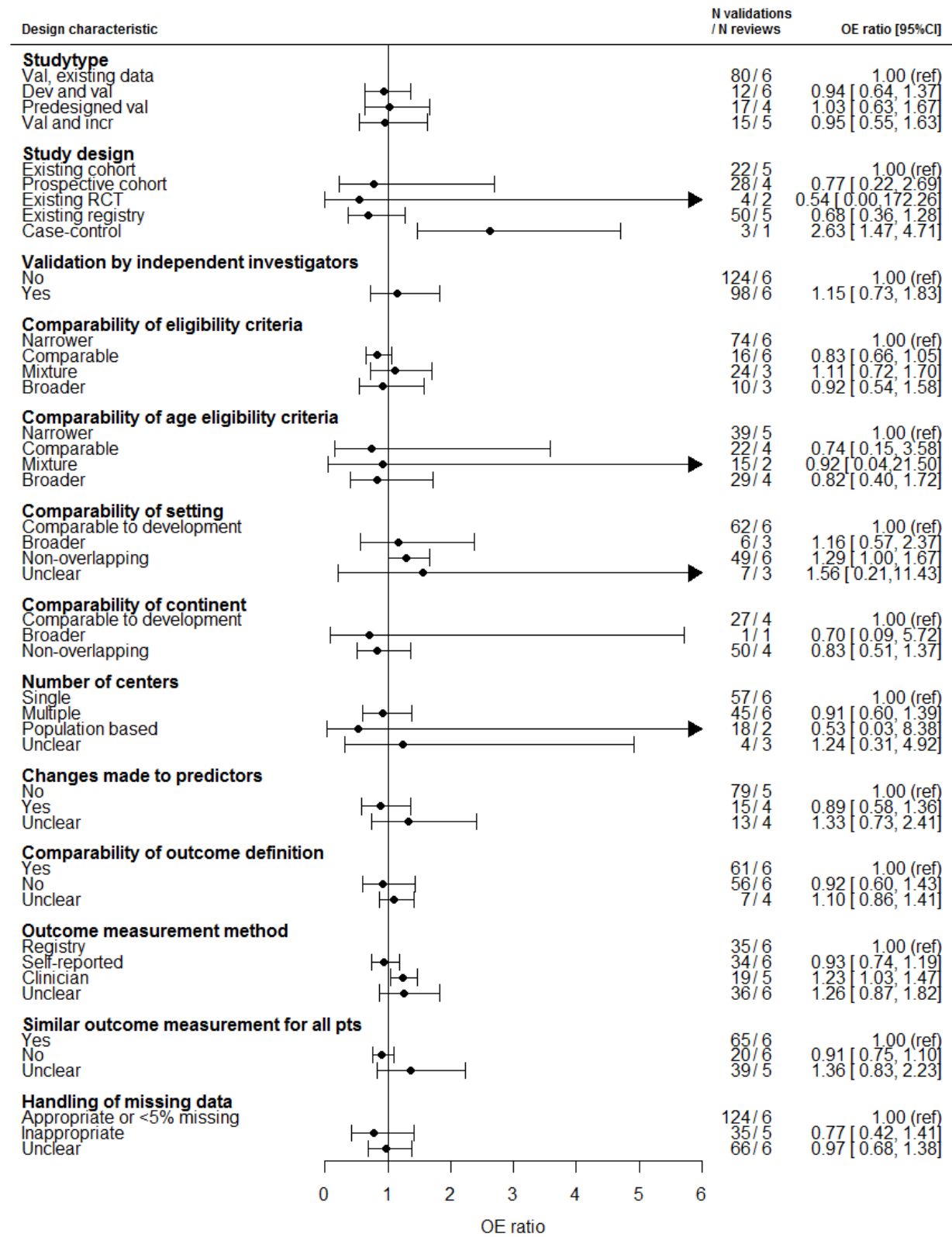
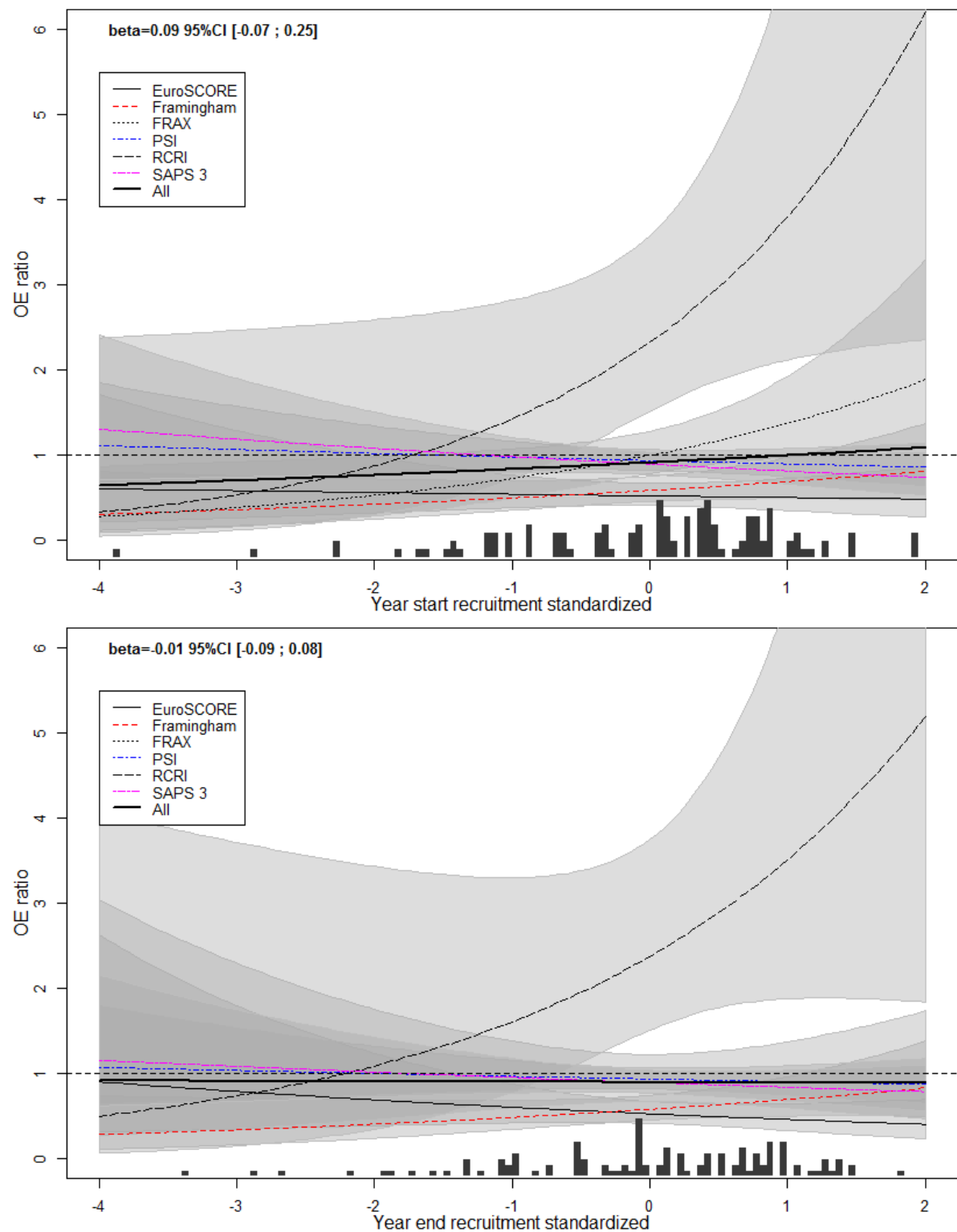
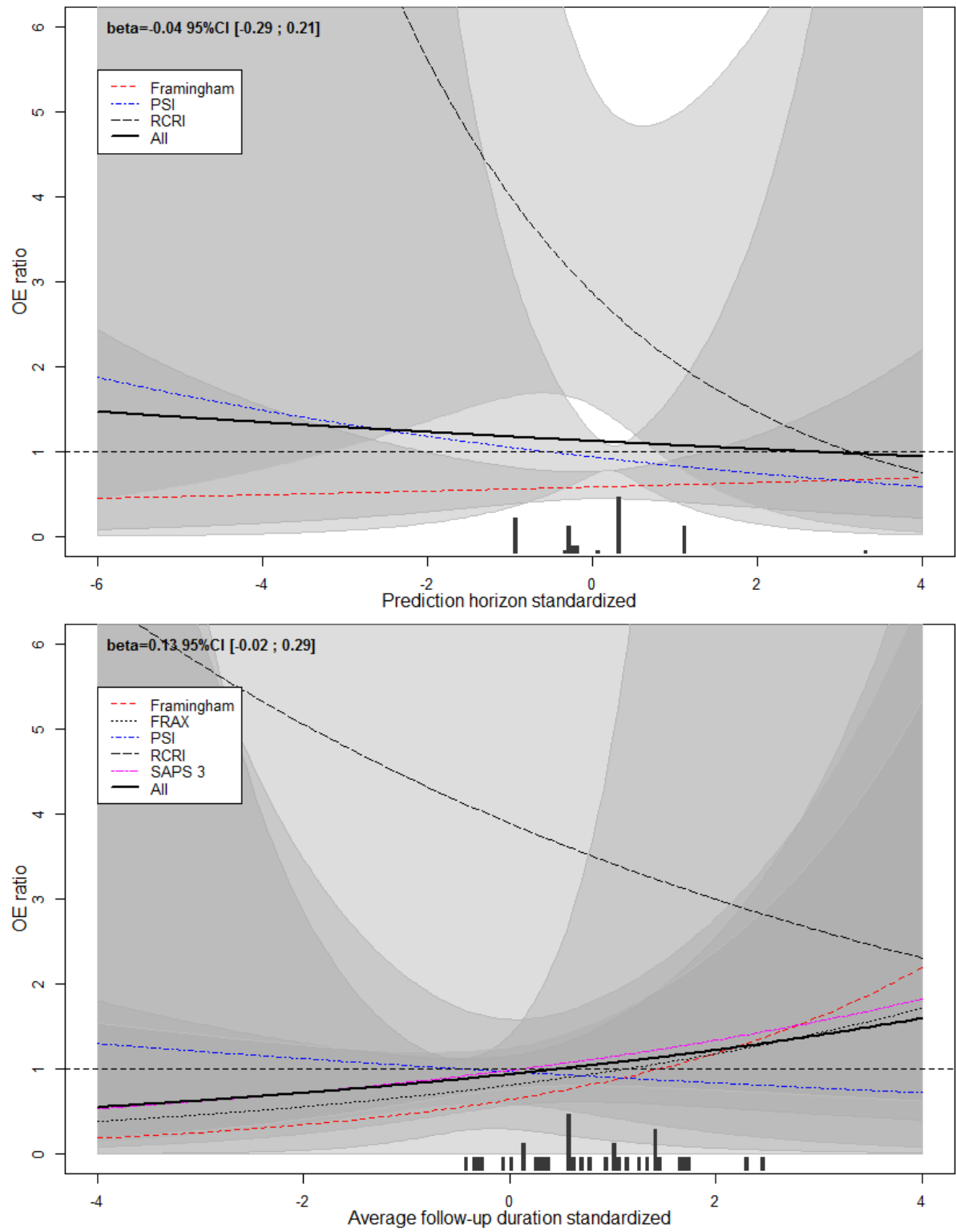
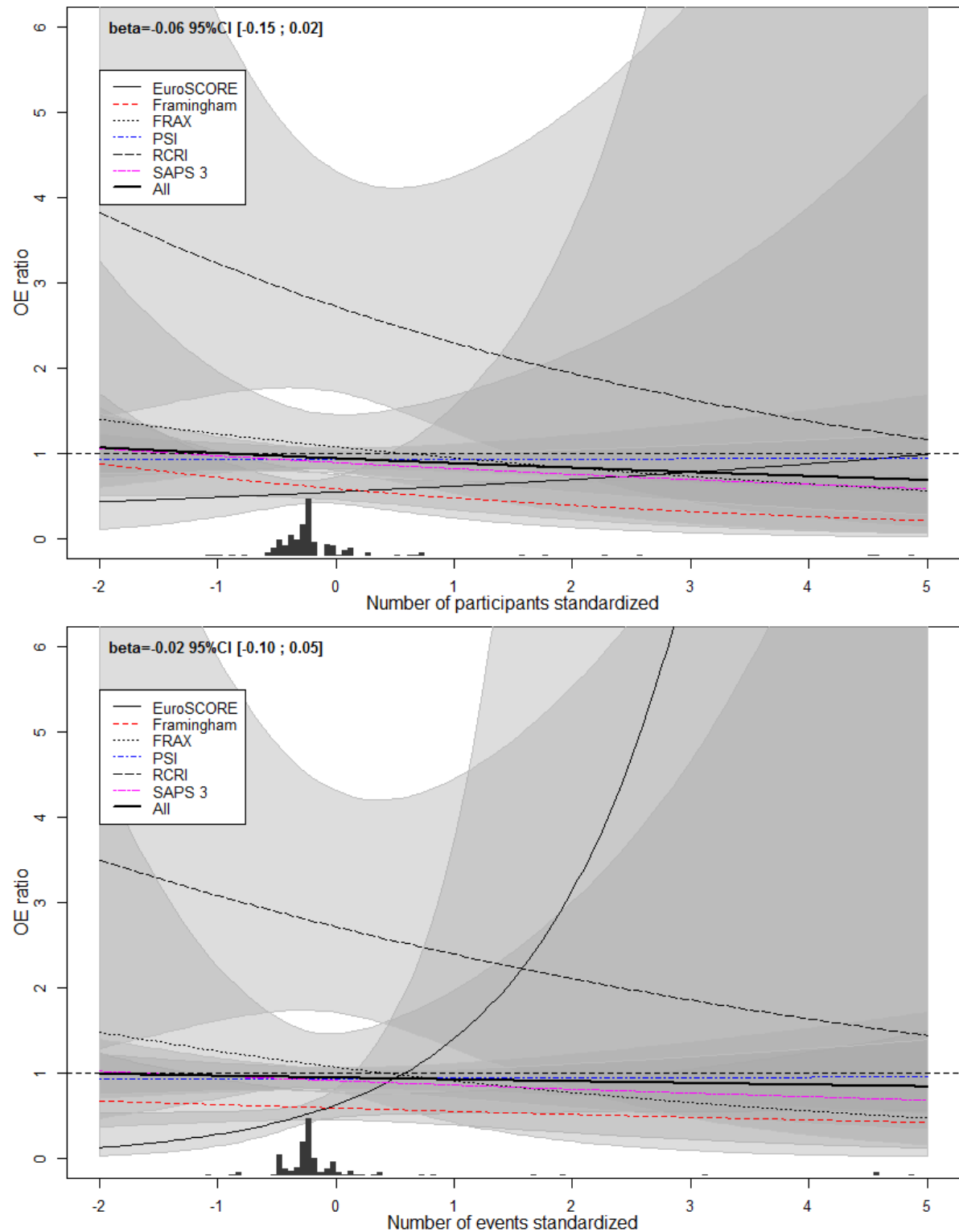
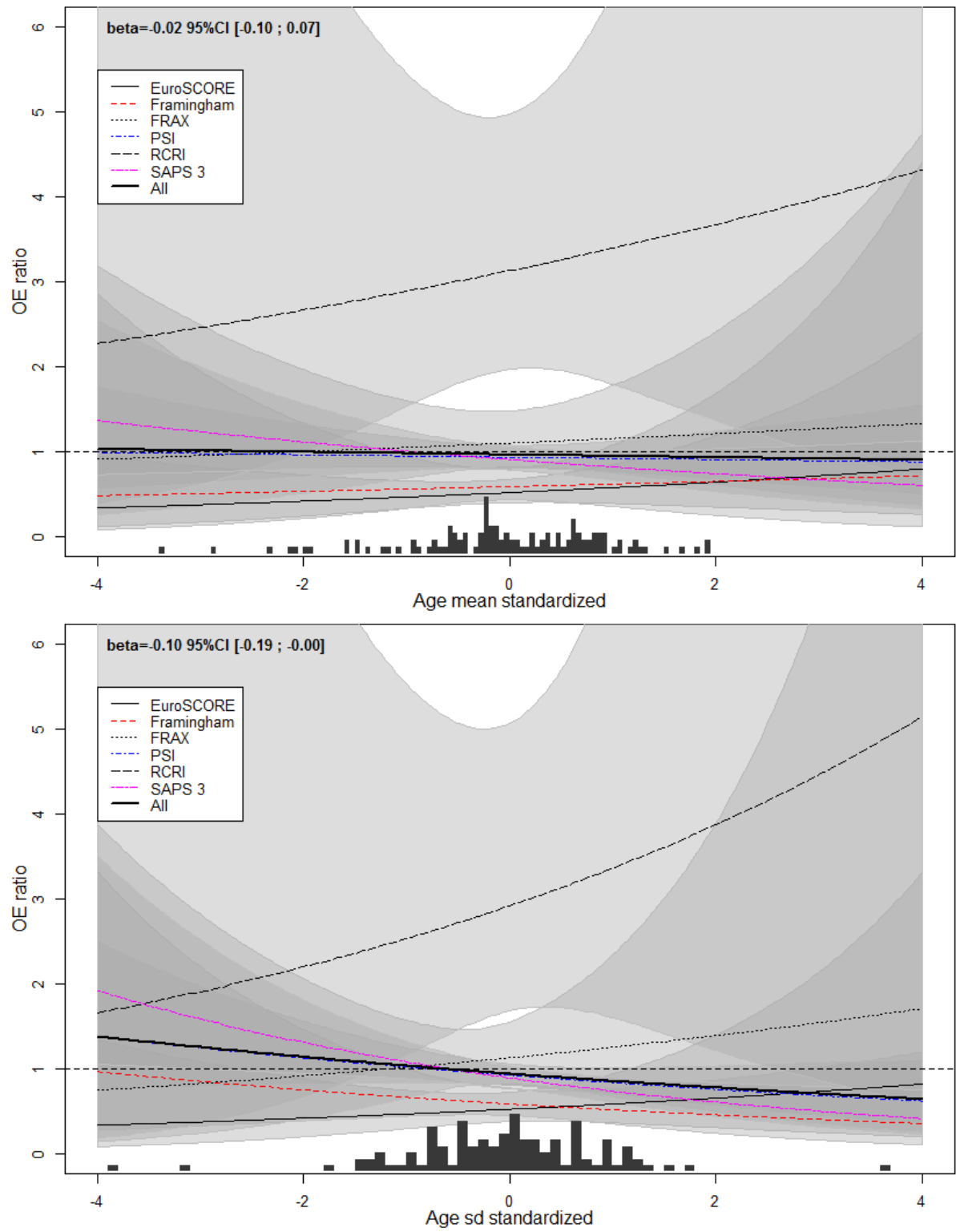


Figure S5: Associations between continuous variables and OE ratio









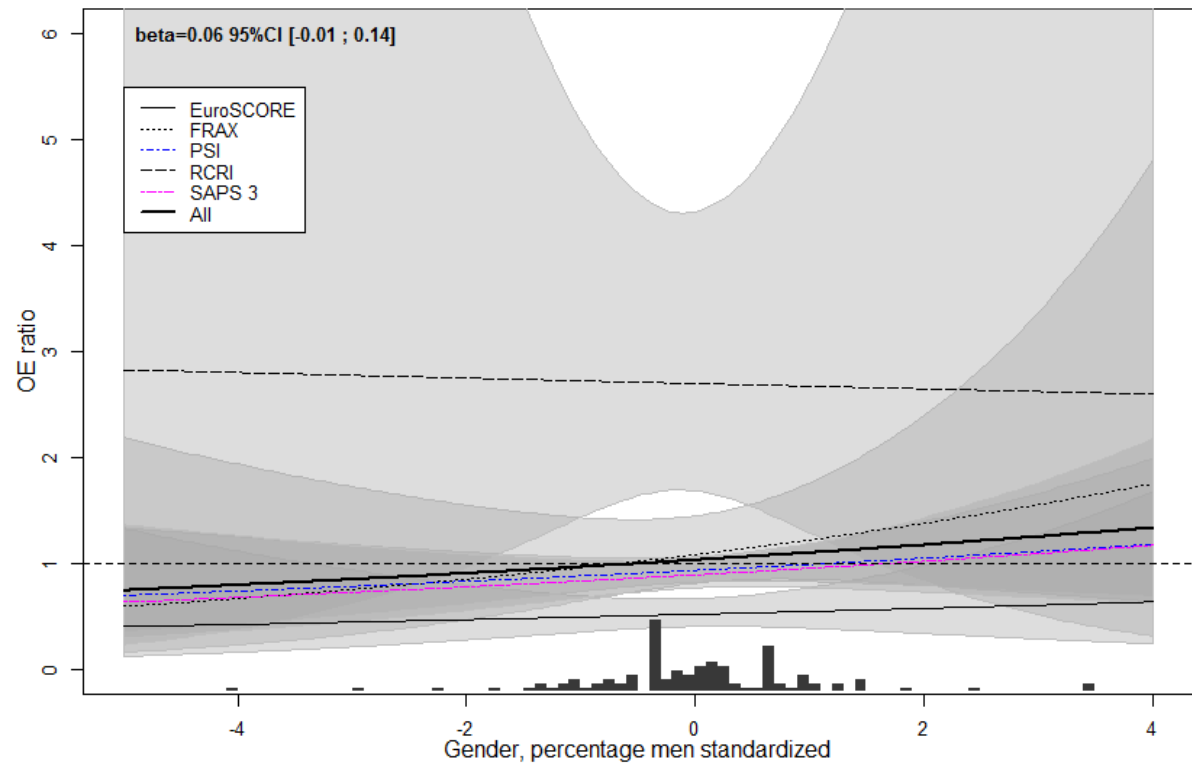
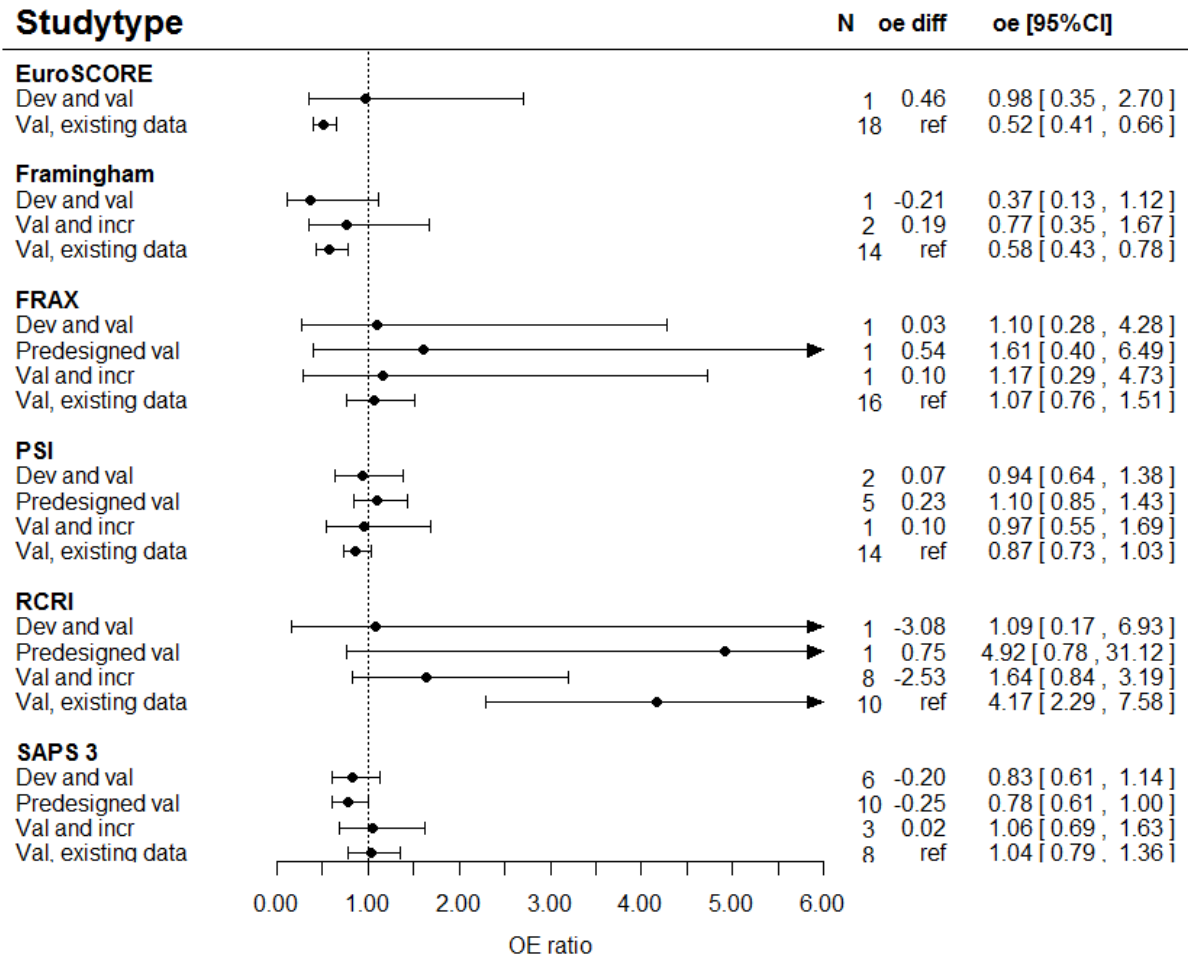
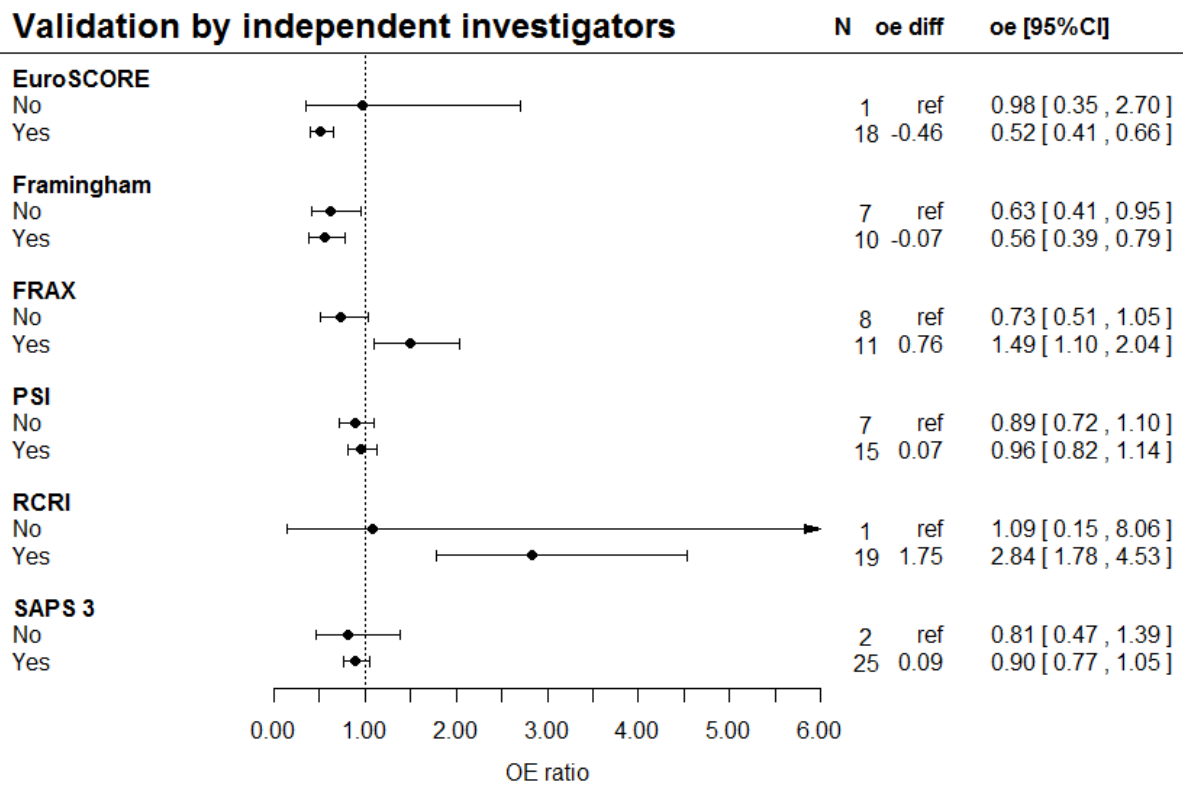
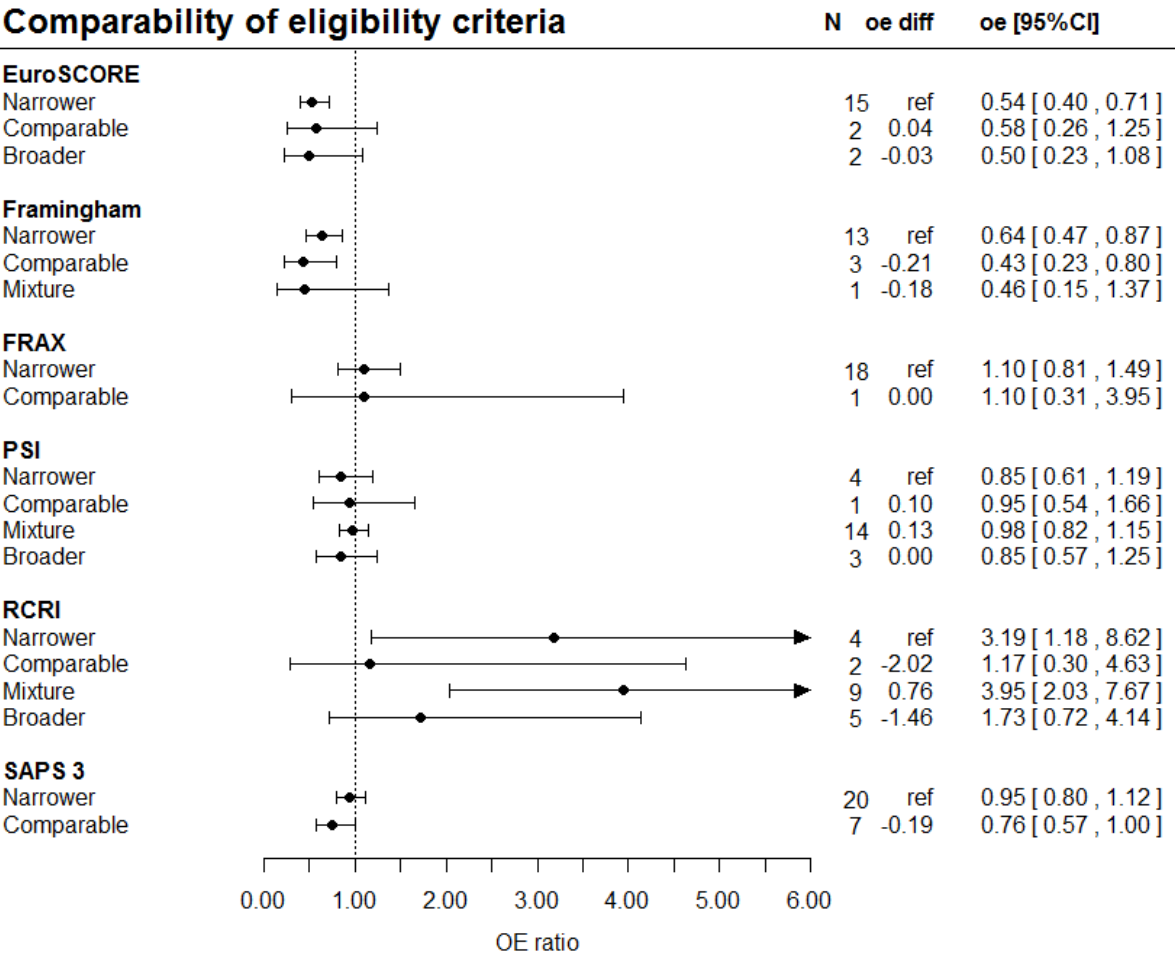


Figure S6: OE ratio in categories of study characteristics within each systematic review

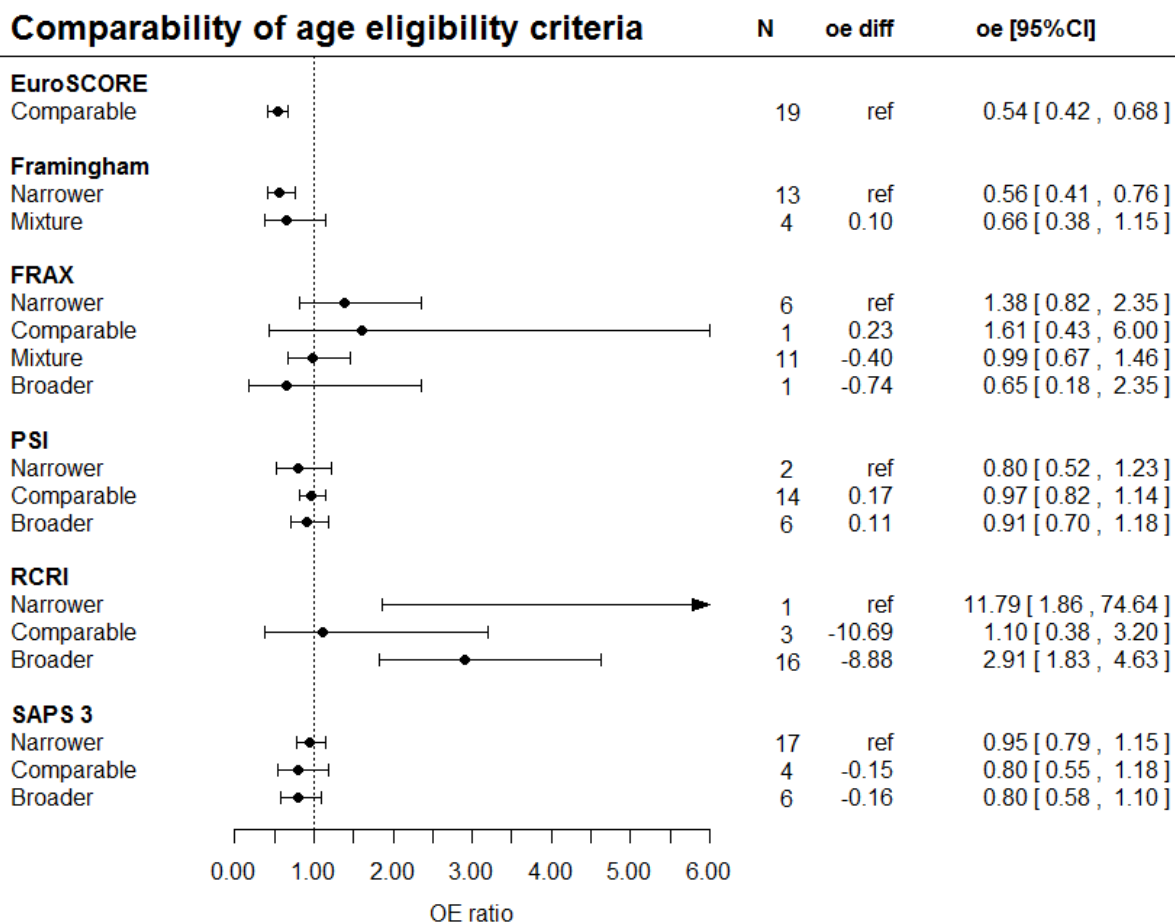


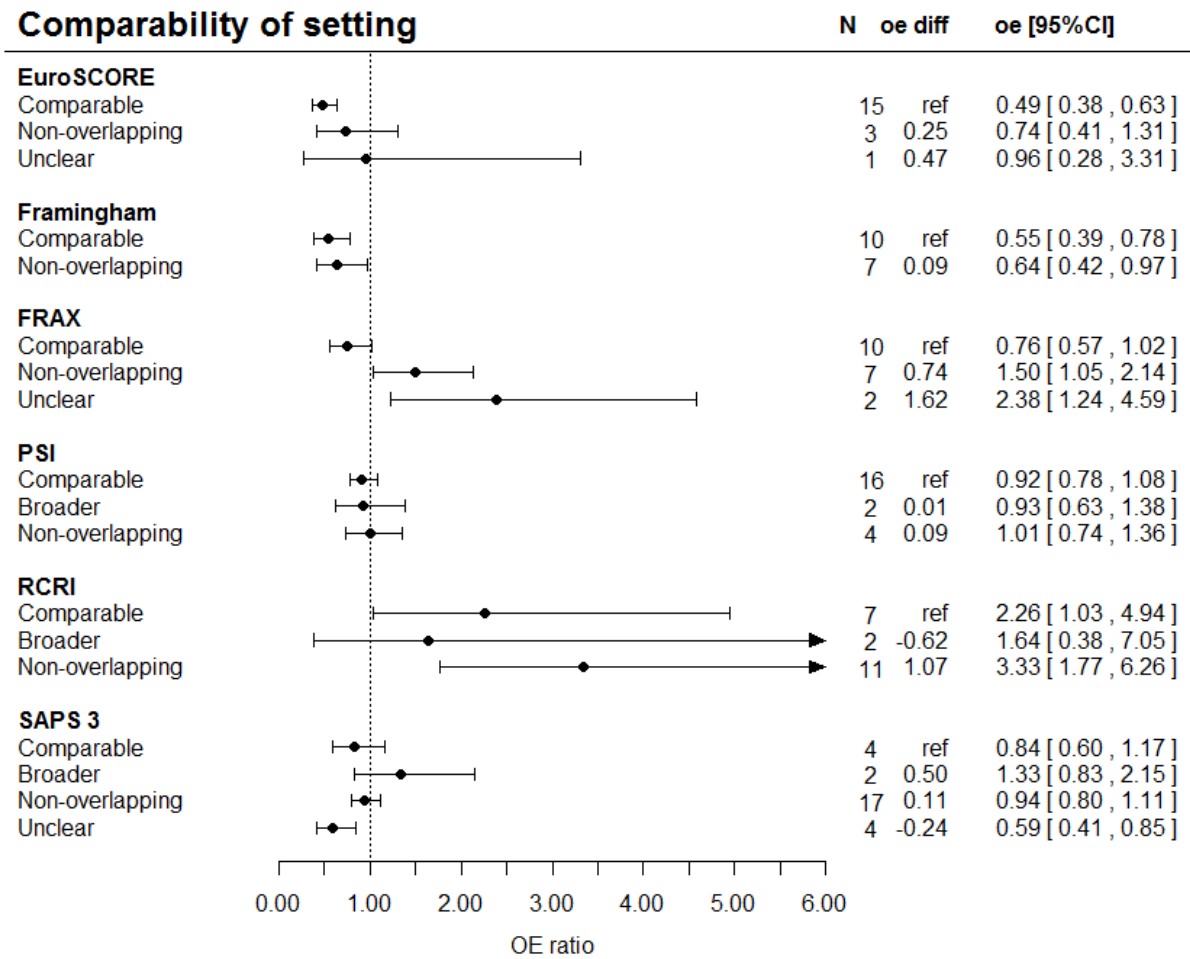
Validation by independent investigators



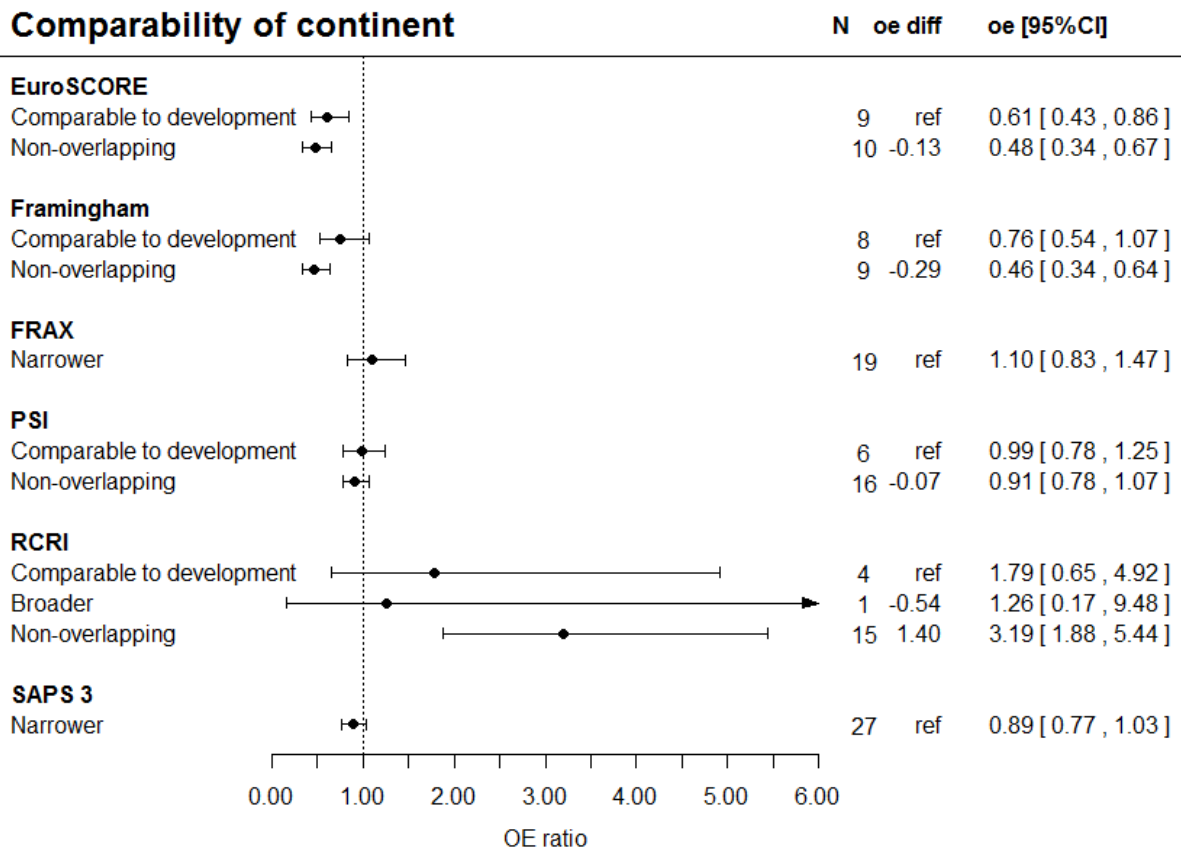


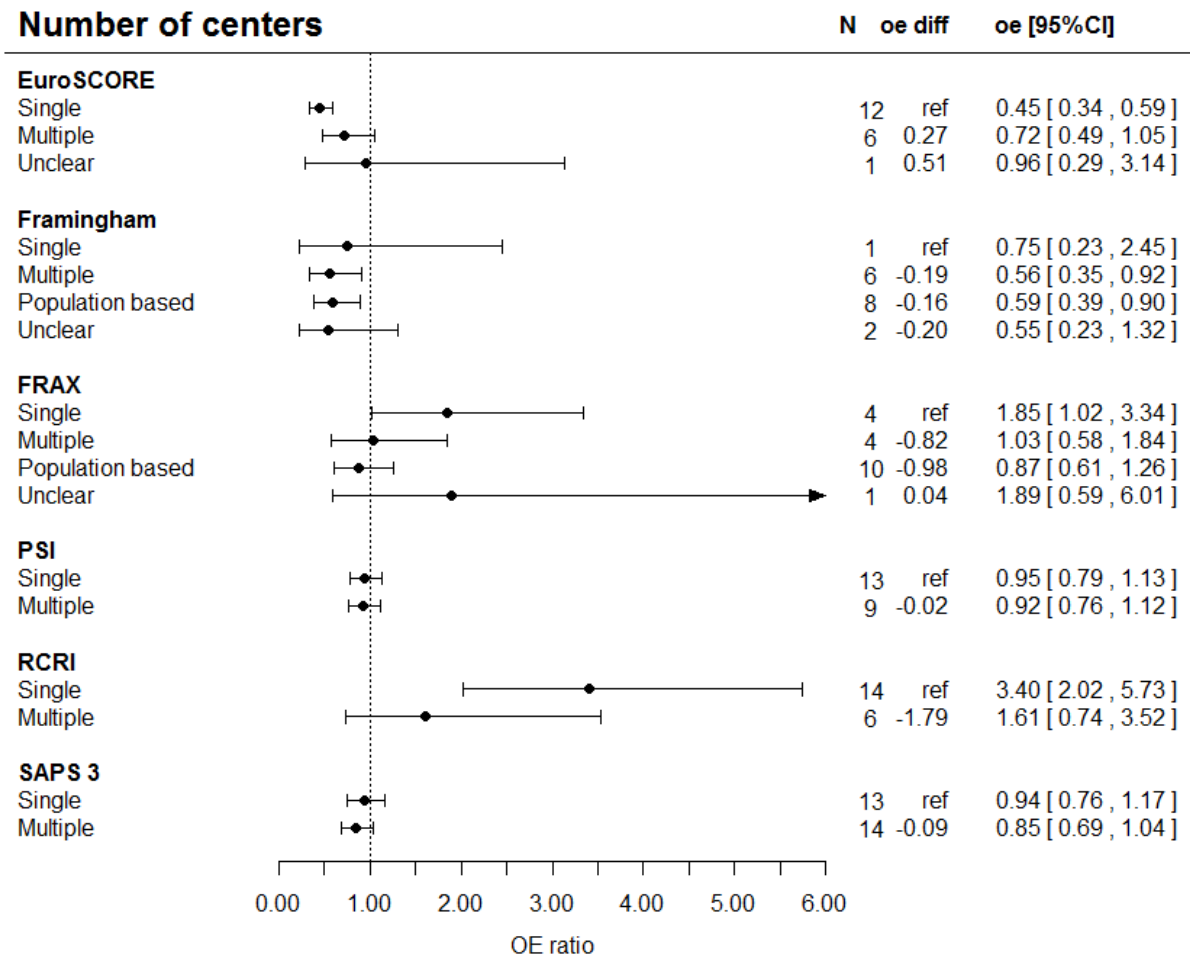
Comparability of age eligibility criteria



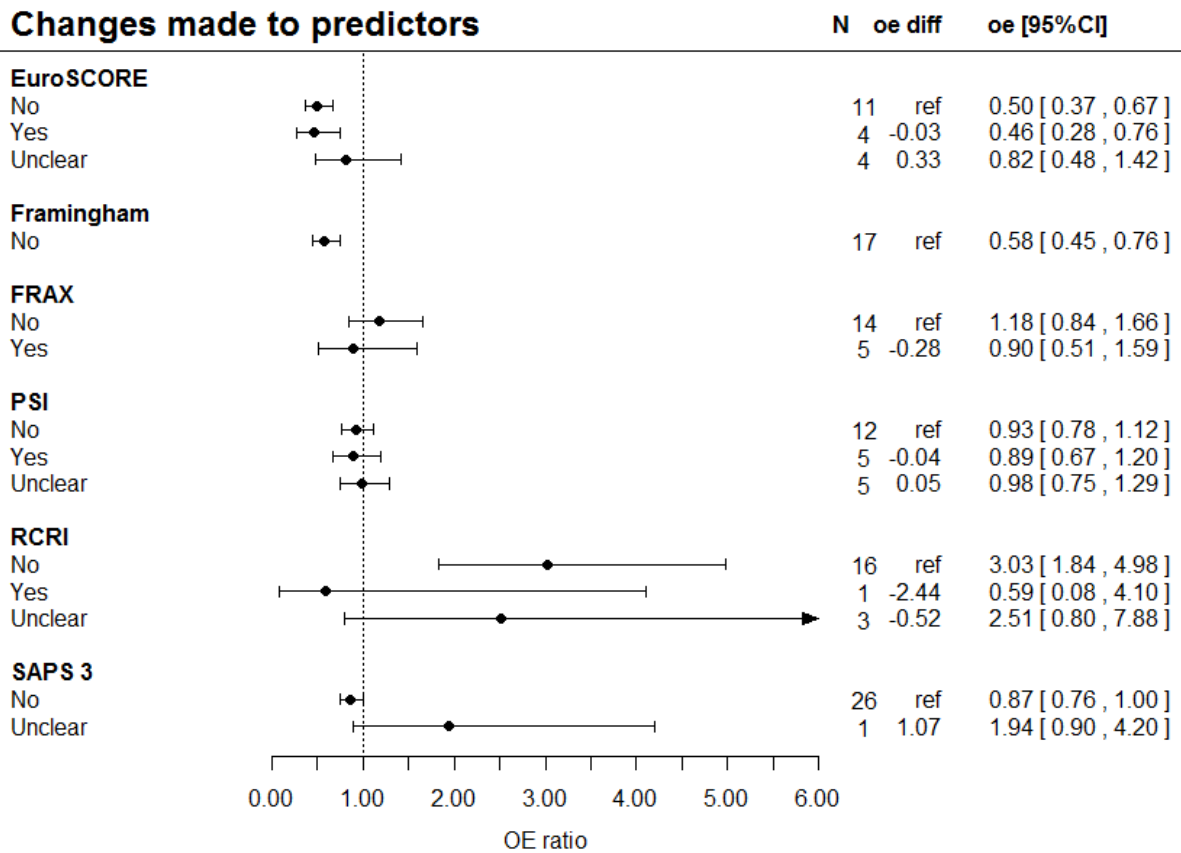


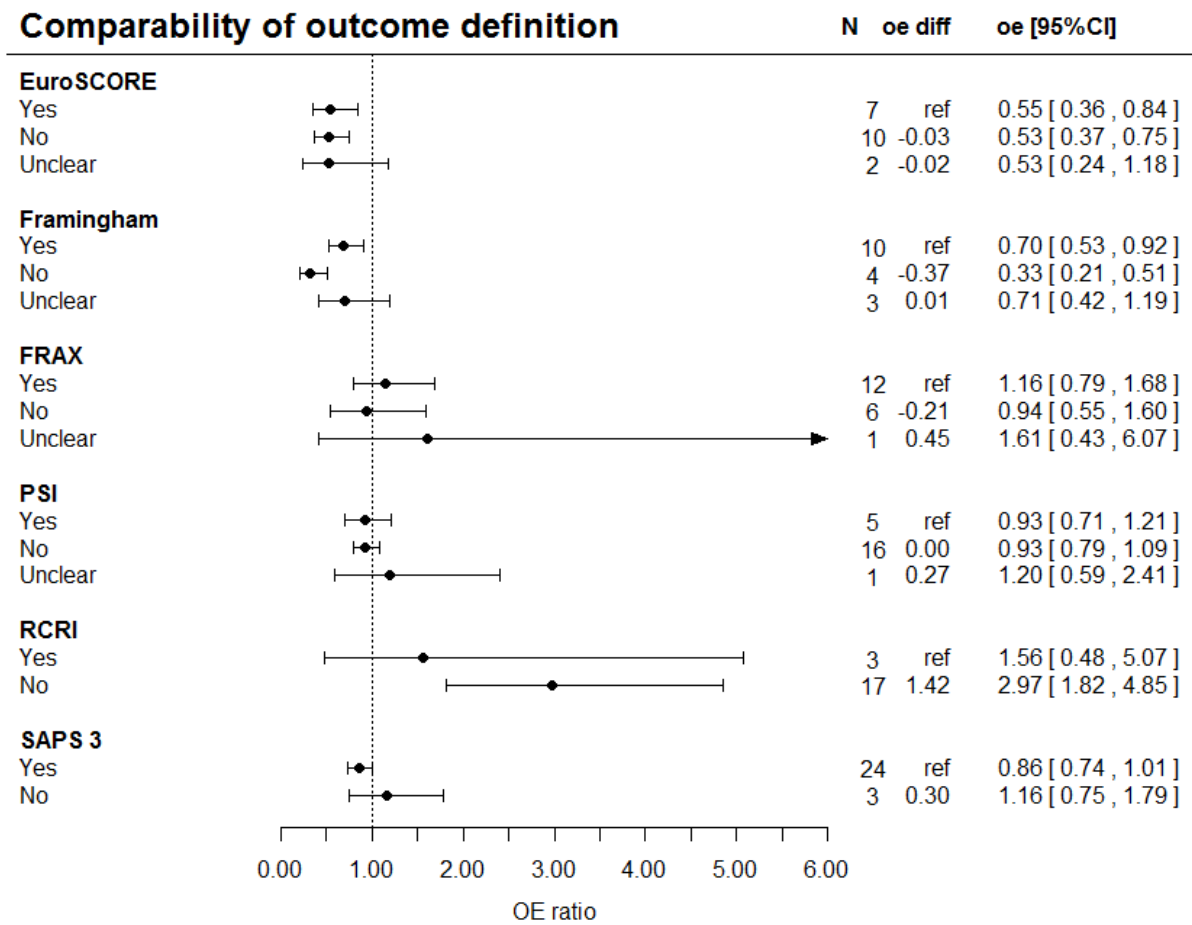
Comparability of continent



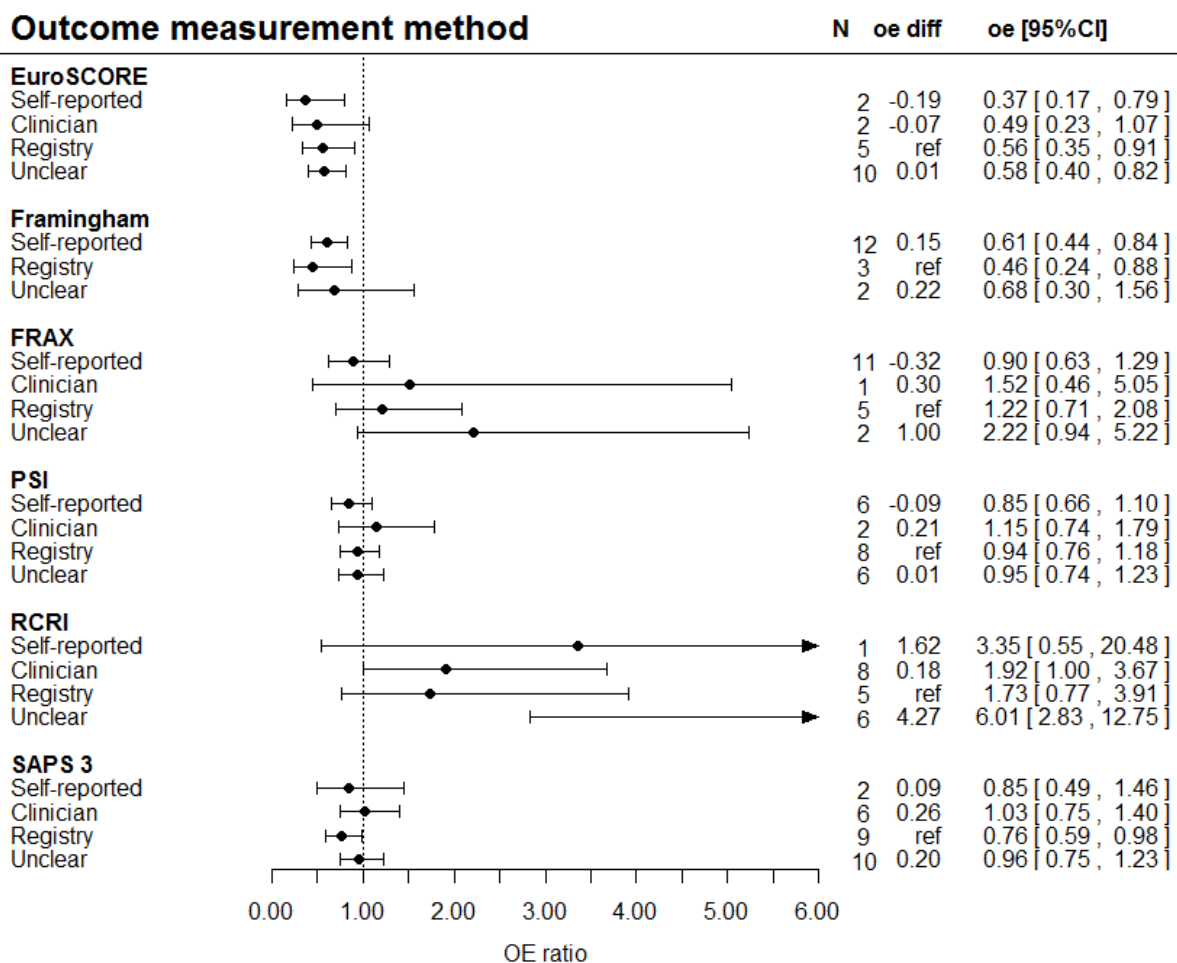


Changes made to predictors

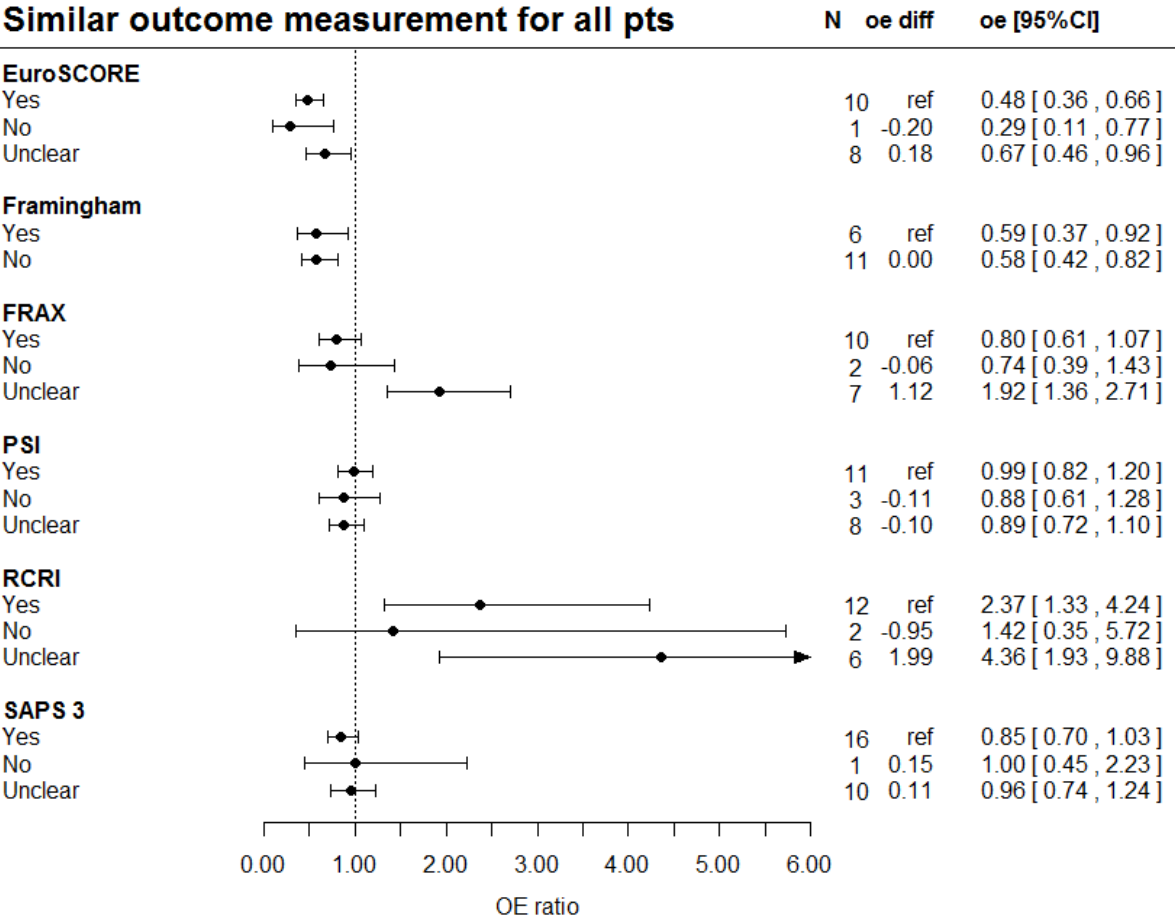


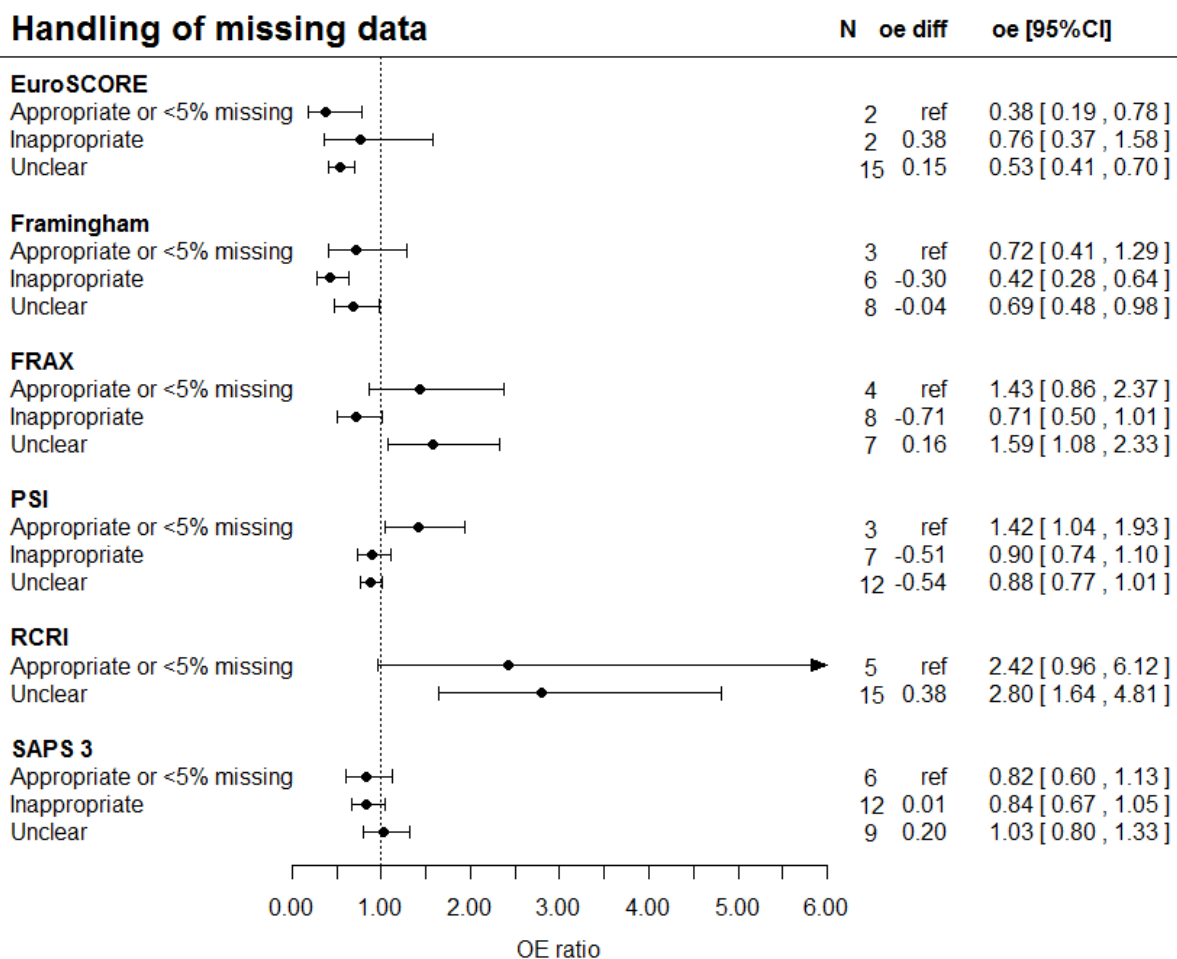


Outcome measurement method



Similar outcome measurement for all pts





OE ratio for categories of study characteristics, pooled using univariable meta-regression analyses per systematic review. N represents the number of external validation studies in a specific category. OE diff represents the difference in OE ratio with regard to a reference category (indicated with 'ref'). Dev: development, val: validation, incr: incremental value, pts: patients.

References

1. Snell KI, Ensor J, Debray TP, Moons KG, Riley RD. Meta-analysis of prediction model performance across multiple studies: Which scale helps ensure between-study normality for the C-statistic and calibration measures? *Stat Methods Med Res* 2017;962280217705678.
2. Newcombe RG. Confidence intervals for an effect size measure based on the Mann-Whitney statistic. Part 2: asymptotic methods and evaluation. *Stat Med* 2006;25(4):559-73.
3. Hanley JA, McNeil BJ. The meaning and use of the area under a receiver operating characteristic (ROC) curve. *Radiology* 1982;143(1):29-36.
4. Debray TP, Damen JA, Snell KI, Ensor J, Hooft L, Reitsma JB, et al. A guide to systematic review and meta-analysis of prediction model performance. *BMJ* 2017;356:i6460.
5. IntHout J, Ioannidis JP, Borm GF. The Hartung-Knapp-Sidik-Jonkman method for random effects meta-analysis is straightforward and considerably outperforms the standard DerSimonian-Laird method. *BMC Med Res Methodol* 2014;14:25.
6. Snell KI, Hua H, Debray TP, Ensor J, Look MP, Moons KG, et al. Multivariate meta-analysis of individual participant data helped externally validate the performance and implementation of a prediction model. *J Clin Epidemiol* 2015.
7. R: A language and environment for statistical computing [program]. Vienna, Austria: R Foundation for Statistical Computing, 2016.
8. Viechtbauer W. Conducting meta-analyses in R with the metafor package. *J Stat Softw* 2010;36(3):1-48.
9. Gasparrini A, Armstrong B, Kenward MG. Multivariate meta-analysis for non-linear and other multi-parameter associations. *Stat Med* 2012;31(29):3821-39.
10. Debray TP. Metamisc: Diagnostic and Prognostic Meta-Analysis. 2017.
11. Bates D, Mächler M, Bolker B, Walker S. Fitting linear mixed-effects models using lme4. *J Stat Softw* 2015;67(1).

Section/topic	Proposed item to be used in methodology research	Reported on page
TITLE		
Title	Identify the report as a meta-epidemiologic study.	1
ABSTRACT		
Structured summary	Provide a structured summary that includes the background of the topic, goal of the study, data sources, method of data selection, appraisal and synthesis methods, results, limitations, conclusions and implications of key findings.	3
INTRODUCTION		
Rationale	Describe the rationale for the meta-epidemiological study in the context of what is already known.	5
Objectives	Provide an explicit statement of the goal of the meta-epidemiological study and the hypothesis being empirically tested.	5
METHODS		
Protocol	Indicate if a protocol exists, if and where it can be accessed (eg, Web address). Registration of a protocol is not mandatory.	Available on request
Eligibility criteria	Specify study characteristics used as criteria for eligibility with a rationale.	7
Information sources	Describe all information sources (eg, databases with dates of coverage, contact with experts to identify additional studies, Internet searches) and search date.	7
Search	Present full electronic search strategy for at least one database, including any limits used, such that it could be repeated. Search is commonly not driven by a clinical question.	Supplement page 2
Study selection	Describe the process for selecting studies for inclusion (ie, how many reviewers selected studies, reviewing in duplicate or by single individuals).	7,8
Data collection process	Describe method of data extraction from reports (eg, piloted forms, independently, in duplicate) and any processes used for manipulating data or obtaining and confirming data from investigators.	8,9, Supplement page 3-6
Data items	List and define all variables for which data were sought and any assumptions and imputations made.	Supplement page 3-6
Risk of bias in individual studies	If risk of bias assessment of individual studies was relevant to the analysis, describe the items used and how this information is to be used during data synthesis.	Not assessed
Summary measures	State the principal summary measures (eg, ratio of risk ratios, difference in means) and explain its meaning and direction to readers.	9, Supplement page 7, 8
Synthesis of results	Describe the statistical or descriptive methods of synthesis including measures of consistency if relevant. If applicable, describe the development of statistical or simulation modelling based on theoretical background. Describe and justify assumptions and	9, Supplement page 7, 8

From: Murad MH, Wang Z. Guidelines for reporting meta-epidemiological methodology research. *Evid Based Med* 2017;22(4):139-42.

Section/topic	Proposed item to be used in methodology research	Reported on page
	computational approximations. Describe methods of additional analyses (eg, sensitivity or subgroup analyses, meta-regression), if done, indicating which were prespecified.	
RESULTS		
Study selection	Give numbers of studies assessed for eligibility and included in the study, with reasons for exclusions at each stage, ideally with a flow diagram. Present a measure of inter-reviewer agreement (eg, kappa statistic).	10, Figure 1
Study characteristics	For each study, present characteristics for which data were extracted and provide the citations. Clinical characteristics may not always be relevant.	Supplement page 9-12
Risk of bias within studies	If risk of bias assessment of individual studies was used in the meta-epidemiological analysis, report risk of bias indicators of each study to allow replication of findings.	Not assessed
Results of individual studies	Present data elements used in the meta-epidemiological analysis from each study (results of clinical outcomes may not be relevant).	Not done
Synthesis of results	Present results of statistical analysis done, including measures of precision and measures of consistency. Present validity of assumptions and fit of statistical or simulation modelling, if applicable.	11, 12, Figure 2-5, Supplement page 13-49
Additional analysis	Give results of additional analyses, if done (eg, sensitivity or subgroup analyses, meta-regression).	Not done
DISCUSSION		
Summary of evidence	Summarise the main findings and compare them with existing knowledge about the topic. The quality of evidence may not be relevant; however, investigators should describe their certainty in the results to readers.	13,14
Limitations	Discuss limitations at research methodology level (eg, likelihood of reporting or publication bias).	13,14
Conclusions	Provide general interpretation of the results and implications for future research. Provide any plausible impact on clinical practice.	16
FUNDING		
Funding	Describe sources of funding for the methodology research and role of funders.	17

From: Murad MH, Wang Z. Guidelines for reporting meta-epidemiological methodology research. Evid Based Med 2017;22(4):139-42.

BMJ Open

Empirical evidence on the impact of study characteristics and the performance of prediction models: a meta-epidemiological study

Journal:	<i>BMJ Open</i>
Manuscript ID	bmjopen-2018-026160.R1
Article Type:	Research
Date Submitted by the Author:	05-Nov-2018
Complete List of Authors:	Damen, Johanna; Cochrane Netherlands, University Medical Center Utrecht, Utrecht University; Julius Center for Health Sciences and Primary Care, University Medical Center Utrecht, Utrecht University Debray, Thomas; Julius Center for Health Sciences and Primary Care, University Medical Center Utrecht, Utrecht University; Cochrane Netherlands, University Medical Center Utrecht, Utrecht University Pajouheshnia, Romin; Julius Center for Health Sciences and Primary Care, University Medical Center Utrecht, Utrecht University Reitsma, Johannes; Julius Center for Health Sciences and Primary Care, University Medical Center Utrecht, Utrecht University; Cochrane Netherlands, University Medical Center Utrecht, Utrecht University Scholten, Rob; Cochrane Netherlands, University Medical Center Utrecht, Utrecht University; Julius Center for Health Sciences and Primary Care, University Medical Center Utrecht, Utrecht University Moons, Karel; Julius Center for Health Sciences and Primary Care, University Medical Center Utrecht, Utrecht University; Cochrane Netherlands, University Medical Center Utrecht, Utrecht University Hooft, Lotty; Cochrane Netherlands, University Medical Center Utrecht, Utrecht University; Julius Center for Health Sciences and Primary Care, University Medical Center Utrecht, Utrecht University
Primary Subject Heading:	Research methods
Secondary Subject Heading:	Epidemiology
Keywords:	Meta-epidemiology, Prognosis, Prognostic models, Bias, Prediction

SCHOLARONE™
Manuscripts

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

1 Empirical evidence on the impact of study characteristics and the performance of prediction models: a
2 meta-epidemiological study

3
4
5
6
7
8 Johanna A A G Damen, Thomas P A Debray, Romin Pajouheshnia, Johannes B Reitsma, Rob J P M
9 Scholten, Karel G M Moons, Lotty Hooft

10
11
12
13
14
15 Johanna A A G Damen
16 Assistant professor
17
18 Cochrane Netherlands, University Medical Center Utrecht, Utrecht University, Utrecht, The Netherlands;
19 Julius Center for Health Sciences and Primary Care, University Medical Center Utrecht, Utrecht
20 University, 3508 GA Utrecht, The Netherlands
21
22
23
24

25 Thomas P A Debray
26 Assistant professor
27
28 Cochrane Netherlands, University Medical Center Utrecht, Utrecht University, Utrecht, The Netherlands;
29 Julius Center for Health Sciences and Primary Care, University Medical Center Utrecht, Utrecht
30 University, 3508 GA Utrecht, The Netherlands
31
32
33
34

35 Romin Pajouheshnia
36 Postdoctoral researcher
37
38 Julius Center for Health Sciences and Primary Care, University Medical Center Utrecht, Utrecht
39 University, 3508 GA Utrecht, The Netherlands
40
41
42
43

44 Johannes B Reitsma
45 Associate professor
46
47 Cochrane Netherlands, University Medical Center Utrecht, Utrecht University, Utrecht, The Netherlands;
48 Julius Center for Health Sciences and Primary Care, University Medical Center Utrecht, Utrecht
49 University, 3508 GA Utrecht, The Netherlands
50
51
52
53

54 Rob J P M Scholten
55 Professor
56
57
58
59
60

33 Cochrane Netherlands, University Medical Center Utrecht, Utrecht University, Utrecht, The Netherlands;
34 Julius Center for Health Sciences and Primary Care, University Medical Center Utrecht, Utrecht
35 University, 3508 GA Utrecht, The Netherlands

37 Karel G M Moons

38 Professor

39 Cochrane Netherlands, University Medical Center Utrecht, Utrecht University, Utrecht, The Netherlands;
40 Julius Center for Health Sciences and Primary Care, University Medical Center Utrecht, Utrecht
41 University, 3508 GA Utrecht, The Netherlands

43 Lotty Hooft

44 Associate professor

45 Cochrane Netherlands, University Medical Center Utrecht, Utrecht University, Utrecht, The Netherlands;
46 Julius Center for Health Sciences and Primary Care, University Medical Center Utrecht, Utrecht
47 University, 3508 GA Utrecht, The Netherlands

49 Correspondence to:

50 Johanna A A G Damen

51 Julius Center for Health Sciences and Primary Care, University Medical Center Utrecht

52 P.O. Box 85500

53 Str. 6.131

54 3508 GA Utrecht

55 The Netherlands

56 j.a.a.damen@umcutrecht.nl

57 +31 88 75 693 77

58

59 Word count: 4167

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

Objectives: To empirically assess the relation between study characteristics and prognostic model performance in external validation studies of multivariable prognostic models.

Design: Meta-epidemiological study.

Data sources and study selection: On October 16th, 2018, we searched electronic databases for systematic reviews of prognostic models. Reviews from non-overlapping clinical fields were selected if they reported common performance measures (either the concordance (c)-statistic or the ratio of observed over expected number of events (OE ratio)) from ten or more validations of the same prognostic model.

Data extraction and analyses: Study design features, population characteristics, methods of predictor and outcome assessment, and the aforementioned performance measures were extracted from the included external validation studies. Random effects meta-regression was used to quantify the association between the study characteristics and model performance.

Results: We included 10 systematic reviews, describing a total of 224 external validations, of which 221 reported c-statistics and 124 OE ratios. Associations between study characteristics and model performance were heterogeneous across systematic reviews. C-statistics were most associated with variation in population characteristics, outcome definitions and measurement, and predictor substitution. For example, validations with eligibility criteria comparable to the development study were associated with higher c-statistics compared to narrower criteria (difference in logit c-statistic 0.21 [95% CI 0.07, 0.35], similar to an increase from 0.70 to 0.74). Using a case-control design was associated with higher OE ratios, compared to using data from a cohort (difference in log OE ratio 0.97 [95% CI 0.38, 1.55], similar to an increase in OE ratio from 1.00 to 2.63).

Conclusions: Variation in performance of prognostic models across studies is mainly associated with variation in case-mix, study designs, outcome definitions and measurement methods, and predictor substitution. Researchers developing and validating prognostic models should realise the potential influence of these study characteristics on the predictive performance of prognostic models.

Strengths and limitations of this study

- To the best of our knowledge, this is the first meta-epidemiological study focusing on the association of study characteristics with estimates of prognostic model performance.
- We included all ten systematic reviews describing at least ten external validations of the same prognostic model, resulting in 224 external validations.
- We extracted relevant features of design and conduct according to existing checklists on quality assessment (CHARMS) and reporting of prediction model studies (TRIPOD).
- It was not feasible to fit multivariable meta-regression models due to the limited number of available, well-reported, validation studies within the individual reviews, rendering the effective sample size too small for multivariable meta-regression analyses.

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

Introduction

Prediction models, including diagnostic and prognostic models, estimate the probability that an individual has or will develop a certain outcome (e.g. disease or complication). Hereto, they combine multiple predictors into an estimate of an individual's risk.¹ Before using a prediction model in clinical practice it is recommended to validate the performance of the model in a population other than the population in which the model was developed (so called external validation studies).² Such studies assess whether model predictions remain sufficiently accurate across different settings and populations. Obviously, it is important that the methodological quality of external validation studies is good, as otherwise estimates of the prediction model's performance may be biased and thereby lead to misleading conclusions on its generalizability to practice.

Systematic reviews have found that the performance of existing prediction models often varies substantially across external validation studies of those models.³⁻⁵ These differences may not only appear due to random variation (when validation studies are small), but may also arise when model predictions are invalid because the model is applied in very different populations (eg, the association between predictors in the model and the outcome are different) or when design-related characteristics of the validation study (eg, measurement methods or variable definitions) are not well aligned with the original development study.^{2,6}

To provide empirical evidence of the association of study characteristics with prediction model performance, a meta-epidemiological approach can be used. Studies using this approach have shown the influence of study characteristics on the effectiveness of interventions studied in randomized trials and on the accuracy of diagnostic tests.⁷⁻¹² For diagnostic prediction models evidence suggests estimates of performance may be biased in studies with certain study characteristics. One study found a higher diagnostic odds ratio in case-control studies, studies with differential outcome verification (ie, using different outcome assessments across study individuals), and with low sample size.¹³ To date, no meta-epidemiological study has been performed investigating the possible impact of study characteristics on measures of the predictive performance of a prognostic model upon external validation, which is commonly quantified in terms of discrimination and calibration.¹⁴ The aim of this study was to investigate sources of heterogeneity in the predictive performance of prognostic models. A meta-

127 epidemiological approach was used to synthesize evidence from a range of clinical fields. This study can
128 serve as empirical evidence for design and analysis related bias in prognostic model studies.

For peer review only

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

Methods

Search and selection of systematic reviews

We used an existing database (last updated on October 16th, 2018) consisting of studies evaluating multiple existing prediction models, including narrative or systematic reviews of prediction models, or head-to-head comparisons of multiple prediction models validated on a specific dataset (See Supplement 1 for details of the search strategy and selection criteria). To construct this database, references identified by the search were screened for eligibility by one reviewer (GSC) on title, abstract and, if necessary, on full text. Subsequently, the full text of all articles in the database were screened for eligibility to the current project by another reviewer (JAAGD). We selected systematic reviews of prognostic models (ie, diagnostic models were excluded) that included at least ten studies that externally validated the same prognostic model (this number was chosen a priori to facilitate the estimation of study-level parameters such as between-study heterogeneity), and that presented the performance of these models in terms of discrimination (concordance (c)-statistic or area under the receiver operating characteristic (AUC) curve), or calibration (observed expected (OE) ratio). Discrimination is the ability of the model to distinguish between people who will and who will not develop the outcome of interest, while calibration reflects the overall agreement between the total number of observed and predicted ('expected') events.¹⁴ We excluded systematic reviews that selected studies based on specific study characteristics (eg, we excluded systematic reviews that did not include primary studies with a sample size below 100, if we were not able to identify the primary studies that had been excluded for this reason). Furthermore, we excluded reviews of prognostic models in which the weights of predictors in the original model were based on expert opinion rather than on coefficients estimated from a formal statistical approach. If more than one systematic review on the same prognostic model was identified, we included the one with the broadest inclusion criteria (eg, reviews focussing on specific patient populations were not preferred if a review with a broader population was available) or the most recent review (in this order of preference). When multiple prognostic models for the same condition were described in one systematic review which all fulfilled the selection criteria, we included the model with the highest number of external validations.

Selection of the primary external validation studies from the included systematic reviews

From the included systematic reviews we collected the primary studies in which the prognostic models were developed and externally validated. For primary external validation studies for which no measure

of discrimination (c-statistic) or calibration (total OE ratio) was reported in the systematic review, we checked the full text of the primary study, and if performance was not reported, these studies were excluded.

If primary external validation studies described multiple external validations of the same model and if there was no overlap in included participants between these external validations (eg, a model was validated in two different cohorts, or a model was validated in men and women separately), data were extracted for every external validation separately. If a model was validated multiple times on the same population (described in either one or multiple publications), we selected the external validation that was included in the systematic review. If the systematic review included all those external validations, we selected the one in which the study population and predicted outcome most closely resembled the population and outcome of the original model.

Data extraction and preparation

We extracted relevant features of design and conduct according to existing checklists on quality assessment (CHARMS) and reporting of prediction model studies (TRIPOD).¹⁵⁻¹⁷ Information about study characteristics of studies in which the models were developed were extracted from the corresponding development papers. Information about study characteristics of primary external validation studies were first extracted from systematic reviews. This information was subsequently checked using the external validation studies and, if necessary, additional information was extracted by one reviewer (JAAGD or RP). Items we extracted included study type (eg, external validation only, development of a new model and external validation of a model), study design (eg, existing cohort, existing RCT), dependency of investigators (validation by independent investigators or investigators also involved in the development study), eligibility criteria for participant inclusion, setting, location (continent), study dates, number of centres, follow-up time and prediction horizon, age and gender distribution, deletion or substitution of predictors, outcome definition and measurement method, sample size and number of events, handling of missing data, and model performance (see Supplement 2 for details). The data extraction form was piloted on multiple articles by all reviewers (JAAGD, TPAD, LH, KGMM, RP, JBR, RJPMS).

For analysis purposes, some study characteristics had to be categorized or transformed (Supplement 2). For example, eligibility criteria of the validation study as compared to the development study had to be judged and categorized as comparable, narrower (if subgroups included in the development study were excluded from the validation study), broader (if subgroups excluded from the development study were

included in the validation study), mixture (a combination of the two), or unclear. For setting, location, predictors and outcome a similar categorization was used. If data on study characteristics were not reported in the primary external validation studies, these were either categorized as ‘unclear’ (in case of categorical study variables), or the study was excluded from the analyses of that (missing) study characteristic (in case of continuous study variables, such as sample size). In order to improve comparability between reviews, we standardized continuous study variables separately for each systematic review, i.e. for every variable we subtracted the mean and divided by the standard deviation of all external validations identified from the same systematic review.

Statistical analyses

We used a two-staged approach to study the possible association between study characteristics and predictive performance.

In the first stage, we fitted a univariable meta-regression model for every study characteristic within each systematic review with the logit c-statistic or log OE ratio as outcome variable.¹⁸ The regression coefficients estimated from this meta-regression model indicate the difference in logit c-statistic or log OE ratio between a certain category of a study characteristic and a chosen reference category (ie, the category that was present in most systematic reviews) of that characteristic.

In the second stage, these regression coefficients were pooled by the use of a random effects model. This reflected the average influence of the study characteristic on model performance across all systematic reviews. For continuous characteristics, the regression coefficients obtained in the first stage were jointly pooled across reviews, using bivariate meta-analysis.^{19 20} For categorical characteristics the results of univariable meta-analyses are presented. We planned to perform multivariable analyses to assess the association between various study characteristics in combination and the performance of prognostic models, but due to the paucity of data we were not able to do so. All analyses are described in more detail in Supplement 3.

Patient and public involvement

Patients and public were not involved in the design, recruitment or conduct of the study.

Results

Identification and selection of studies

The search identified 2392 studies, of which 555 were included in the database and screened on full text, and 79 were further assessed (Figure 1). Finally, ten systematic reviews were included.²¹⁻²⁹ These reviews addressed external validations of the following prognostic models: ABCD2,³⁰ Essen Stroke Risk Score (ESRS),³¹ EuroSCORE,³² Framingham,³³ FRAX,³⁴ Injury Severity Score (ISS),³⁵ model for end-stage liver disease (MELD),³⁶ Pneumonia Severity Index (PSI),³⁷ Revised Cardiac Risk Index (RCRI),³⁸ and Simplified Acute Physiology Score (SAPS) 3³⁹ (Table 1). The reviews included 248 primary external validation studies with 274 external model validations (one study could describe multiple model validations). During data extraction, 73 of 274 validations were eventually excluded (most often for not reporting a performance measure), and 20 additional external model validations were identified (Figure 1). This resulted in the inclusion of 224 external validations, of which 221 could be included in the analyses of the c-statistic, and 124 in the analyses of the total OE ratio. For the total OE ratio, only validations of the EuroSCORE, Framingham, FRAX, PSI, RCRI and SAPS 3 prognostic models were included, due to the very low number of reported OE ratios in the validations studies for the other four prognostic models.

Description of included validations

The number of external validations per systematic review ranged from 11 to 30 (Table 1), and the median (IQR) sample size and number of events were 1069 (418-3043) and 92 (36-248), respectively. Most studies used an existing registry (N=104, 46%) or existing cohort (N=74, 33%) to validate the prognostic model. The median (IQR) c-statistic and total OE ratio were 0.73 (0.64-0.82) and 0.92 (0.64-1.26), respectively. Predictive performance of the models was highly heterogeneous, even for external validations of the same prognostic model, as indicated by the wide prediction intervals (Table 1). Not all information on the study characteristics was reported for all external validations (Table S1). Information was often unclear (eg, for outcome definitions (N=83, 37%) and handling of missing data (N=105, 47%)) or missing (eg, case-mix information such as mean age (N=28, 13%) and gender distribution (N=16, 7%)).

Discrimination

Pooled models

1
2
3 253 The pooled analyses across all systematic reviews (Figure 2, S1 and S2) showed that validation in a
4
5 254 continent different from the development study was associated with a higher c-statistic, compared to
6
7 255 validation in the same continent, and multicentre versus single centre validation studies were associated
8
9 256 with a lower c-statistic. Comparable eligibility criteria for participant inclusion were also associated with
10
11 257 higher c-statistics compared to narrower criteria, whereas a broader setting was associated with a lower
12
13 258 c-statistic compared to a setting comparable to the development study. Although not statistically
14
15 259 significant, validations with changes made to the predictors (ie, substitution or deletion of a predictor),
16
17 260 or in which it was unclear whether all predictors were correctly measured, tended to have lower c-
18
19 261 statistics compared to validations where no changes were made. In various reviews we found an
20
21 262 association between the c-statistic and numerous other study characteristics, such as the study design,
22
23 263 comparability of outcome definition, prediction horizon, sample size and number of events, and mean
24
25 264 age of study participants (Figure 3, S2 and S3), only these were often not statistically significant when
26
27 265 pooled together.

266
267 *Variation across reviews*

268 Across reviews we found associations of many study characteristics with the c-statistic although this was
269 rather heterogeneous, and confidence intervals often overlapped (Figure 3 and Figure S3). For example,
270 for study design, in six systematic reviews a higher c-statistic was found for validations that used an
271 existing registry compared to an existing cohort, while in three reviews a lower c-statistic was found. In
272 three systematic reviews we found a higher c-statistic in validations by independent investigators, while
273 in five a lower c-statistic was found.

274
275 For other study characteristics, directions of associations were more consistent. For example, for most
276 systematic reviews, validation studies with eligibility criteria narrower compared to the criteria used in
277 the development study had a lower c-statistics while broader eligibility criteria were associated with
278 higher c-statistics (Figure S3). C-statistics were also (slightly) higher in external validations with a setting
279 comparable to the development study. Validation in a continent other than the development study in
280 general was associated with a higher c-statistic, and multicentre studies had lower c-statistics compared
281 to single centre studies. External validations in which it was unclear if there were changes made to the
282 predictors had lower c-statistics (Figure S3).

283
284 Calibration

285 *Pooled analyses*

286 We found a significant association between study design and the total OE ratio (Figure 4); using data
287 from a case-control study (although known to be an inferior design for prognostic model research^{1 6})
288 resulted in higher OE ratios, compared to using data from an existing cohort (though based on three
289 external validations). Furthermore, higher OE ratios were found for studies in which the outcome was
290 assessed by a panel of clinicians as compared to using a registry. In various reviews we found an
291 association between the total OE ratio and numerous other study characteristics, such as the duration of
292 follow-up, year in which recruitment was started, sample size, standard deviation of age, and setting
293 (Figure 4, S4, S5 and S6), only these were not statistically significant when pooled together.

295 *Variation across reviews*

296 For other categories of study design (other than the use of a case-control design), heterogeneous
297 associations were found across systematic reviews (Figure 5). The associations of most other study
298 characteristics with the OE ratio were also most often not consistent across systematic reviews (Figure
299 S5 and S6). For example, for two systematic reviews external validations with appropriate handling of
300 missing data had OE ratios closer to 1 compared to inappropriate handling of missing data, while in two
301 reviews, OE ratios were further away from 1. Only for the continent in which the model was validated,
302 directions were more consistent; OE ratios were closer to 1 if the continent was comparable to the
303 development, compared to validations in different continents (Figure S6).

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

Discussion

Principal findings

Using a comprehensive meta-analytical approach, we studied the association between study characteristics of prognostic model validation studies and the estimated model performance across ten clinical domains. We focused on objective study characteristics that can be extracted from published reports. The reporting of the primary external validation studies was often incomplete and inadequate. Key study characteristics, such as outcome definitions, handling of missing data, and even model calibration estimates were infrequently reported. Still, we found associations between various study characteristics and a model’s predictive performance. Changes in a model’s predictive performance were notably found in relation to validation studies with a case-control (versus cohort) design, with differences in case-mix, in continent (in which the model is validated), in eligibility criteria, in clinical setting, in number of centres (included in the validation study), in differences in outcome definitions and assessments, and in predictor substitutions.

Explanations, strengths and weaknesses

Based on findings in meta-epidemiological studies on the effect of study characteristics and the efficacy of interventions⁷⁻¹⁰ and diagnostic test accuracy,^{11 12} we anticipated to find more statistically significant associations between study characteristics and model performance across the included systematic reviews from different domains. Although we included every systematic review that described at least ten external validation studies of the same prognostic model, resulting in more than 200 validations from 10 reviews, our analyses appeared to still be hampered by relatively low numbers of external validations per systematic review, combined with poor reporting and substantial heterogeneity within and across systematic reviews. Conceptually, there are many potential sources of heterogeneity in model performance, such as differences in population characteristics, predictor and outcome definitions and measurements, and in many aspects of the statistical analyses (eg, dealing with missing data, sample size and selective loss to follow up). All these characteristics may act in isolation but could also be related to each other. The individual strength of the association of one characteristic with model performance is ideally addressed by adopting multivariable (meta)-regression models with the observed model performance estimates of the validation studies as dependent variable and the characteristics of multiple design features as independent variables.^{10 12} Unfortunately, this approach was not feasible here due to the limited

number of available, well-reported, validation studies within the individual reviews, rendering the effective sample size too small for multivariable meta-regression analyses.

A general limitation of all meta-epidemiological studies, is the possibility that the effect of a certain study characteristic differs across systematic reviews which may nullify the effect when pooled together.⁴⁰ We also found numerous conflicting associations between a study characteristic and the reported predictive performance measures across reviews that were cancelled out in the pooled analyses.

Also, it is possible that the effect caused by individual study characteristics is small and therefore difficult to detect. Moreover, there might be some misclassification of study characteristics, caused either by our misinterpretation of what is reported, or by a lack of reporting, which could have diluted the effects of the study characteristics. Indeed, the c-statistic is often considered to be an insensitive measure to quantify changes in model performance.⁴¹⁻⁴³ In previous simulation studies, the c-statistic and OE ratio appeared to be strongly influenced by case-mix differences,^{14 44 45} which may mask the possible (smaller) effects from design-related characteristics. Other measures that are less sensitive to case-mix differences, such as the calibration slope, could, however, not be studied here simply because they were (almost) never reported in our retrieved studies, as was also shown previously.³

We found greater variation in the methods used by external validation studies *between* models than within validations of the *same* model. For example, multiple imputation is the preferred method for handling missing data in prediction modelling.^{46 47} However, in the field of cardiovascular disease, it seems common to handle missing data by performing a complete case analysis, while in the field of mortality prediction in surgical patients, typically researchers fill in 'normal' values if a value is missing. Finally, given the explorative nature of our analyses to identify potential areas of further research, we did not correct for multiple testing, though we tried to minimize the number of exploratory analyses.

Comparison to previous research

Despite above considerations, our findings, ie, the trends in the associations between study characteristics and model performance measures (though not always statistically significant), are in agreement with various previous simulation studies in this field.^{14 44 46-48} For example, we confirmed that studies with more variation in case-mix show higher c-statistics, and lower c-statistics when a predictor was omitted from the model. However, we found lower c-statistics in studies with a broader setting and when the number of centres in a study was higher. The lower c-statistic in multicentre studies might be caused by increased variation in predictor definitions and methods to measure predictors compared to

1
2
3 368 single centre studies, where it is more likely that definitions and measurement methods have been
4
5 369 standardized. This might result in increased measurement error in multicentre studies, which is known
6
7 370 to lower discriminative ability of a model.^{49 50}
8
9 371 We also found a higher total OE ratio in studies with a case-control design. Both simulation studies and
10
11 372 meta-epidemiological studies in the fields of diagnostic tests and (mainly diagnostic) prediction models,
12
13 373 have shown biased effect measures in studies using a case-control design.¹¹⁻¹⁴ This confirms that case-
14
15 374 control studies should not be used to study certain aspects of model calibration. Further, we found that
16
17 375 the total OE ratio was influenced by the method of outcome assessment, in agreement with previous
18
19 376 studies that showed that higher diagnostic odds ratios were found in studies with differential outcome
20
21 377 verification.¹³ We also expected to find lower OE ratios when the validation population differed from the
22
23 378 development population (eg, in terms of case-mix).¹⁴ We could not systematically confirm this across all
24
25 379 reviews, likely caused by heterogeneity between systematic reviews as indicated by the wide confidence
26
27 380 intervals. Finally, we could not fully confirm the association between sample size and model
28
29 381 performance that was previously found,¹³ although we found similar trends in part of the reviews.
30
31 382

32
33 383 Implications for future research

34
35 384 In agreement with many previously conducted systematic reviews on prediction models,^{3 51-55} we still
36
37 385 and again found poor reporting of prediction model studies. Meta-epidemiological studies of prediction
38
39 386 model studies would highly benefit from complete reporting according to the Transparent reporting of a
40
41 387 multivariable prediction model for individual prognosis or diagnosis (TRIPOD) statement.^{16 17} We
42
43 388 recommend at least the following information, essential for comparing and interpreting the results of
44
45 389 external validation studies, to be reported by every external validation study: eligibility criteria for
46
47 390 participant inclusion, details of predictor and outcome definitions and measurements, a clear reference
48
49 391 to the model that is being validated and any changes made to this validated model compared to the
50
51 392 model as presented in the development study, estimates of model discrimination and calibration
52
53 393 performance (including calibration slope and intercept), and corresponding standard errors for the
54
55 394 original model and, if applicable, for any updated model.
56
57 395 We also believe that more research is urgently needed to evaluate under which circumstances certain
58
59 396 design choices may lead to heterogeneity in prediction model performance, and to incorporate these
60
397 issues in the appraisal of prediction model studies. There is a need for more guidance on how to score
398 items of critical appraisal checklists for prediction model studies, such as the CHARMS checklist.¹⁵

Several options exist to gain more empirical insight in design related bias in prediction model studies. Firstly, meta-epidemiological researchers can collect more external validation studies and try to correct for all issues that cause variation in performance of a model. We believe, however, that this is currently not feasible as we already included every systematic reviews describing at least ten validations of the same prognostic model. A second and much more efficient option is to collect the individual participant data (IPD) for all studies included in this review to directly study the effect of study characteristics on model performance.⁵⁶⁻⁶⁰ Using IPD, it will also be possible to study different performance measures, like the case-mix adjusted c-statistic^{44 61} and calibration slope.¹⁴ Thirdly, new simulation studies could be performed to get more insight in design related bias in prediction model performance. Researchers could for example study the effect of using a different outcome definition or prediction horizon on the c-statistic of a model.

Conclusion

In this comprehensive meta-epidemiological study we found empirical evidence for an association between study characteristics and predictive performance of prognostic models. We found that predictive performance of prognostic models upon external validation is highly heterogeneous, but sensitive to various study characteristics, such as study design, case-mix, eligibility criteria, setting, methods of outcome definition and measurement, and predictor substitution. It is important that these characteristics are thus emphasized in the reporting and appraisal of prediction model studies. However, for a large part the observed heterogeneity in model performance remained unexplained, which is likely caused by the high number of factors that cause heterogeneity in predictive performance and may act in opposite directions whereas a multivariable meta-regression analysis across reviews simply was not possible.

1
2
3 422 Acknowledgments: The authors would like to acknowledge Prof. Gary S Collins (GSC) for building the
4
5 423 database with systematic reviews of prediction models, which served as a basis for this paper.
6
7 424
8 425 Contributors: KGMM, JBR, TPAD, and LH conceived the study. JAAGD, TPAD, RP, JBR, RJPMS, KGMM and
9
10 426 LH were involved in designing the study. JAAGD selected the articles. JAAGD and RP extracted the data.
11
12 427 JAAGD analysed the data in close consultation with TPAD. JAAGD, TPAD, RP, JBR, RJPMS, KGMM and LH
13
14 428 were involved in interpreting the data. JAAGD wrote the first draft of the manuscript which was revised
15
16 429 by TPAD, RP, JBR, RJPMS, KGMM and LH. The corresponding author attests that all listed authors meet
17
18 430 authorship criteria and that no others meeting the criteria have been omitted. JAAGD is guarantor.
19
20 431
21 432 Funding: Thomas Debray gratefully acknowledges the Netherlands Organization for Health Research and
22
23 433 Development (grant number 91617050). Karel GM Moons received a grant from The Netherlands
24
25 434 Organization for Scientific Research (ZONMW 918.10.615 and 91208004). The funder had no role in the
26
27 435 design of the study; the collection, analysis, and interpretation of the data; or approval of the finished
28
29 436 manuscript.
30
31 437
32 438 Competing interests: All authors have completed the ICMJE uniform disclosure form at
33
34 439 www.icmje.org/coi_disclosure.pdf and declare: no support from any organisation for the submitted
35
36 440 work; no financial relationships with any organisations that might have an interest in the submitted
37
38 441 work in the previous three years; no other relationships or activities that could appear to have
39
40 442 influenced the submitted work.
41
42 443
43 444 Ethical approval: Not required.
44
45 445
46 446 Data sharing: no additional data are available.
47
48 447
49 448 Transparency: The lead author (JAAGD) affirms that the manuscript is an honest, accurate, and
50
51 449 transparent account of the study being reported; that no important aspects of the study have been
52
53 450 omitted; and that any discrepancies from the study as planned (and, if relevant, registered) have been
54
55 451 explained.
56
57 452
58
59
60

1
2
3 453 The Corresponding Author has the right to grant on behalf of all authors and does grant on behalf of all
4
5 454 authors, a worldwide licence to the Publishers and its licensees in perpetuity, in all forms, formats and
6
7 455 media (whether known now or created in the future), to i) publish, reproduce, distribute, display and
8
9 456 store the Contribution, ii) translate the Contribution into other languages, create adaptations, reprints,
10
11 457 include within collections and create summaries, extracts and/or, abstracts of the Contribution, iii)
12
13 458 create any other derivative work(s) based on the Contribution, iv) to exploit all subsidiary rights in the
14
15 459 Contribution, v) the inclusion of electronic links from the Contribution to third party material where-
16
17 460 ever it may be located; and, vi) licence any third party to do any or all of the above.
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

References

1. Moons KG, Kengne AP, Woodward M, Royston P, Vergouwe Y, Altman DG, et al. Risk prediction models: I. Development, internal validation, and assessing the incremental value of a new (bio)marker. *Heart* 2012;98(9):683-90.

2. Debray TP, Vergouwe Y, Koffijberg H, Nieboer D, Steyerberg EW, Moons KG. A new framework to enhance the interpretation of external validation studies of clinical prediction models. *J Clin Epidemiol* 2014.

3. Damen JA, Hooft L, Schuit E, Debray TP, Collins GS, Tzoulaki I, et al. Prediction models for cardiovascular disease risk in the general population: systematic review. *BMJ* 2016;353:i2416.

4. Alba AC, Agoritsas T, Jankowski M, Courvoisier D, Walter SD, Guyatt GH, et al. Risk prediction models for mortality in ambulatory patients with heart failure: a systematic review. *Circ Heart Fail* 2013;6(5):881-9.

5. Meads C, Ahmed I, Riley RD. A systematic review of breast cancer incidence risk prediction models with meta-analysis of their performance. *Breast Cancer Res Treat* 2012;132(2):365-77.

6. Moons KG, Royston P, Vergouwe Y, Grobbee DE, Altman DG. Prognosis and prognostic research: what, why, and how? *BMJ* 2009;338:b375.

7. Page MJ, Higgins JP, Clayton G, Sterne JA, Hrobjartsson A, Savovic J. Empirical Evidence of Study Design Biases in Randomized Trials: Systematic Review of Meta-Epidemiological Studies. *PLoS One* 2016;11(7):e0159267.

8. Berkman ND, Santaguida PL, Viswanathan M, Morton SC. AHRQ Methods for Effective Health Care. *The Empirical Evidence of Bias in Trials Measuring Treatment Differences*. Rockville (MD): Agency for Healthcare Research and Quality (US), 2014.

9. Savovic J, Jones H, Altman D, Harris R, Juni P, Pildal J, et al. Influence of reported study design characteristics on intervention effect estimates from randomised controlled trials: combined analysis of meta-epidemiological studies. *Health Technol Assess* 2012;16(35):1-82.

10. Wood L, Egger M, Gluud LL, Schulz KF, Juni P, Altman DG, et al. Empirical evidence of bias in treatment effect estimates in controlled trials with different interventions and outcomes: meta-epidemiological study. *BMJ* 2008;336(7644):601-5.

11. Lijmer JG, Mol BW, Heisterkamp S, Bossel GJ, Prins MH, van der Meulen JH, et al. Empirical evidence of design-related bias in studies of diagnostic tests. *JAMA* 1999;282(11):1061-6.

12. Rutjes AW, Reitsma JB, Di Nisio M, Smidt N, van Rijn JC, Bossuyt PM. Evidence of bias and variation in diagnostic accuracy studies. *CMAJ* 2006;174(4):469-76.

13. Ban JW, Emparanza JI, Urreta I, Burls A. Design Characteristics Influence Performance of Clinical Prediction Rules in Validation: A Meta-Epidemiological Study. *PLoS One* 2016;11(1):e0145779.

14. Steyerberg E. *Clinical prediction models: a practical approach to development, validation, and updating*: Springer Science & Business Media, 2008.

15. Moons KG, de Groot JA, Bouwmeester W, Vergouwe Y, Mallett S, Altman DG, et al. Critical appraisal and data extraction for systematic reviews of prediction modelling studies: the CHARMS checklist. *PLoS Med* 2014;11(10):e1001744.

16. Collins GS, Reitsma JB, Altman DG, Moons KG. Transparent Reporting of a multivariable prediction model for Individual Prognosis or Diagnosis (TRIPOD): the TRIPOD statement. *Ann Intern Med* 2015;162(1):55-63.

17. Moons KG, Altman DG, Reitsma JB, Ioannidis JP, 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. Snell KI, Ensor J, Debray TP, Moons KG, Riley RD. Meta-analysis of prediction model performance across multiple studies: Which scale helps ensure between-study normality for the C-statistic and calibration measures? *Stat Methods Med Res* 2017;962280217705678.
19. Snell KI, Hua H, Debray TP, Ensor J, Look MP, Moons KG, et al. Multivariate meta-analysis of individual participant data helped externally validate the performance and implementation of a prediction model. *J Clin Epidemiol* 2015.
20. Debray TP, Damen JA, Snell KI, Ensor J, Hooft L, Reitsma JB, et al. A guide to systematic review and meta-analysis of prediction model performance. *BMJ* 2017;356:i6460.
21. Thompson DD, Murray GD, Dennis M, Sudlow CL, Whiteley WN. Formal and informal prediction of recurrent stroke and myocardial infarction after stroke: a systematic review and evaluation of clinical prediction models in a new cohort. *BMC Med* 2014;12:58.
22. Siregar S, Groenwold RH, de Heer F, Bots ML, van der Graaf Y, van Herwerden LA. Performance of the original EuroSCORE. *Eur J Cardiothorac Surg* 2012;41(4):746-54.
23. Marques A, Ferreira RJ, Santos E, Loza E, Carmona L, da Silva JA. The accuracy of osteoporotic fracture risk prediction tools: a systematic review and meta-analysis. *Ann Rheum Dis* 2015;74(11):1958-67.
24. Tohira H, Jacobs I, Mountain D, Gibson N, Yeo A. Systematic review of predictive performance of injury severity scoring tools. *Scand J Trauma Resusc Emerg Med* 2012;20:63.
25. Klein KB, Stafinski TD, Menon D. Predicting survival after liver transplantation based on pre-transplant MELD score: a systematic review of the literature. *PLoS One* 2013;8(12):e80661.
26. Chalmers JD, Mandal P, Singanayagam A, Akram AR, Choudhury G, Short PM, et al. Severity assessment tools to guide ICU admission in community-acquired pneumonia: systematic review and meta-analysis. *Intensive Care Med* 2011;37(9):1409-20.
27. Ford MK, Beattie WS, Wijeyesundera DN. Systematic review: prediction of perioperative cardiac complications and mortality by the revised cardiac risk index. *Ann Intern Med* 2010;152(1):26-35.
28. Nassar AP, Malbouisson LM, Moreno R. Evaluation of Simplified Acute Physiology Score 3 performance: a systematic review of external validation studies. *Crit Care* 2014;18(3):R117.
29. Damen JAAG, Pajouheshnia R, Heus P, Moons KGM, Reitsma JB, Scholten RJPM, et al. Performance of the Framingham risk models and Pooled Cohort Equations for predicting 10-year risk of cardiovascular disease: a systematic review and meta-analysis. Manuscript submitted for publication.
30. Rothwell PM, Giles MF, Flossmann E, Lovelock CE, Redgrave JN, Warlow CP, et al. A simple score (ABCD) to identify individuals at high early risk of stroke after transient ischaemic attack. *Lancet* 2005;366(9479):29-36.
31. Diener HC, Ringleb PA, Savi P. Clopidogrel for the secondary prevention of stroke. *Expert Opin Pharmacother* 2005;6(5):755-64.
32. Nashef SA, Roques F, Michel P, Gauducheau E, Lemeshow S, Salamon R. European system for cardiac operative risk evaluation (EuroSCORE). *Eur J Cardiothorac Surg* 1999;16(1):9-13.
33. Wilson PW, D'Agostino RB, Levy D, Belanger AM, Silbershatz H, Kannel WB. Prediction of coronary heart disease using risk factor categories. *Circulation* 1998;97(18):1837-47.
34. Kanis JA, Oden A, Johnell O, Johansson H, De Laet C, Brown J, et al. The use of clinical risk factors enhances the performance of BMD in the prediction of hip and osteoporotic fractures in men and women. *Osteoporos Int* 2007;18(8):1033-46.
35. Baker SP, O'Neill B, Haddon W, Jr., Long WB. The injury severity score: a method for describing patients with multiple injuries and evaluating emergency care. *J Trauma* 1974;14(3):187-96.

1
2
3 553 36. Malinchoc M, Kamath PS, Gordon FD, Peine CJ, Rank J, ter Borg PC. A model to predict poor survival
4 554 in patients undergoing transjugular intrahepatic portosystemic shunts. *Hepatology*
5 555 2000;31(4):864-71.
6 556 37. Fine MJ, Auble TE, Yealy DM, Hanusa BH, Weissfeld LA, Singer DE, et al. A prediction rule to identify
7 557 low-risk patients with community-acquired pneumonia. *N Engl J Med* 1997;336(4):243-50.
8 558 38. Lee TH, Marcantonio ER, Mangione CM, Thomas EJ, Polanczyk CA, Cook EF, et al. Derivation and
9 559 prospective validation of a simple index for prediction of cardiac risk of major noncardiac
10 560 surgery. *Circulation* 1999;100(10):1043-9.
11 561 39. Moreno RP, Metnitz PG, Almeida E, Jordan B, Bauer P, Campos RA, et al. SAPS 3--From evaluation of
12 562 the patient to evaluation of the intensive care unit. Part 2: Development of a prognostic model
13 563 for hospital mortality at ICU admission. *Intensive Care Med* 2005;31(10):1345-55.
14 564 40. Sterne JA, Juni P, Schulz KF, Altman DG, Bartlett C, Egger M. Statistical methods for assessing the
15 565 influence of study characteristics on treatment effects in 'meta-epidemiological' research. *Stat*
16 566 *Med* 2002;21(11):1513-24.
17 567 41. Steyerberg EW, Vickers AJ, Cook NR, Gerds T, Gonen M, Obuchowski N, et al. Assessing the
18 568 performance of prediction models: a framework for traditional and novel measures.
19 569 *Epidemiology* 2010;21(1):128-38.
20 570 42. Pencina MJ, D'Agostino RB, Sr., D'Agostino RB, Jr., Vasan RS. Evaluating the added predictive ability
21 571 of a new marker: from area under the ROC curve to reclassification and beyond. *Stat Med*
22 572 2008;27(2):157-72; discussion 207-12.
23 573 43. Uno H, Cai T, Pencina MJ, D'Agostino RB, Wei LJ. On the C-statistics for evaluating overall adequacy
24 574 of risk prediction procedures with censored survival data. *Stat Med* 2011;30(10):1105-17.
25 575 44. Vergouwe Y, Moons KG, Steyerberg EW. External validity of risk models: Use of benchmark values to
26 576 disentangle a case-mix effect from incorrect coefficients. *Am J Epidemiol* 2010;172(8):971-80.
27 577 45. Usher-Smith JA, Sharp SJ, Griffin SJ. The spectrum effect in tests for risk prediction, screening, and
28 578 diagnosis. *BMJ* 2016;353:i3139.
29 579 46. Held U, Kessels A, Garcia Aymerich J, Basagana X, Ter Riet G, Moons KG, et al. Methods for Handling
30 580 Missing Variables in Risk Prediction Models. *Am J Epidemiol* 2016;184(7):545-51.
31 581 47. Janssen KJ, Vergouwe Y, Donders AR, Harrell FE, Jr., Chen Q, Grobbee DE, et al. Dealing with missing
32 582 predictor values when applying clinical prediction models. *Clin Chem* 2009;55(5):994-1001.
33 583 48. Austin PC, Steyerberg EW. Interpreting the concordance statistic of a logistic regression model:
34 584 relation to the variance and odds ratio of a continuous explanatory variable. *BMC Med Res*
35 585 *Methodol* 2012;12:82.
36 586 49. Khudyakov P, Gorfine M, Zucker D, Spiegelman D. The impact of covariate measurement error on
37 587 risk prediction. *Stat Med* 2015;34(15):2353-67.
38 588 50. Pajouheshnia R, van Smeden M, Peelen LM, Groenwold RHH. How variation in predictor
39 589 measurement affects the discriminative ability and transportability of a prediction model. *J Clin*
40 590 *Epidemiol* 2018.
41 591 51. Collins GS, Mallett S, Omar O, Yu LM. Developing risk prediction models for type 2 diabetes: a
42 592 systematic review of methodology and reporting. *BMC Med* 2011;9:103.
43 593 52. Collins GS, Omar O, Shanyinde M, Yu LM. A systematic review finds prediction models for chronic
44 594 kidney disease were poorly reported and often developed using inappropriate methods. *J Clin*
45 595 *Epidemiol* 2013;66(3):268-77.
46 596 53. Bouwmeester W, Zuithoff NP, Mallett S, Geerlings MI, Vergouwe Y, Steyerberg EW, et al. Reporting
47 597 and methods in clinical prediction research: a systematic review. *PLoS Med* 2012;9(5):1-12.
48 598 54. Perel P, Edwards P, Wentz R, Roberts I. Systematic review of prognostic models in traumatic brain
49 599 injury. *BMC Med Inform Decis Mak* 2006;6:38.

55. Mallett S, Royston P, Dutton S, Waters R, Altman DG. Reporting methods in studies developing prognostic models in cancer: a review. *BMC Med* 2010;8:20.
56. Debray TP, Koffijberg H, Vergouwe Y, Moons KG, Steyerberg EW. Aggregating published prediction models with individual participant data: a comparison of different approaches. *Stat Med* 2012;31(23):2697-712.
57. Debray TP, Moons KG, Ahmed I, Koffijberg H, Riley RD. A framework for developing, implementing, and evaluating clinical prediction models in an individual participant data meta-analysis. *Stat Med* 2013;32(18):3158-80.
58. Debray TP, Riley RD, Rovers MM, Reitsma JB, Moons KG. Individual Participant Data (IPD) Meta-analyses of Diagnostic and Prognostic Modeling Studies: Guidance on Their Use. *PLoS Med* 2015;12(10):e1001886.
59. Riley RD, Ensor J, Snell KI, Debray TP, Altman DG, Moons KG, et al. External validation of clinical prediction models using big datasets from e-health records or IPD meta-analysis: opportunities and challenges. *BMJ* 2016;353:i3140.
60. Riley RD, Price MJ, Jackson D, Wardle M, Gueyffier F, Wang J, et al. Multivariate meta-analysis using individual participant data. *Res Synth Methods* 2015;6(2):157-74.
61. White IR, Rapsomaniki E. Covariate-adjusted measures of discrimination for survival data. *Biom J* 2015;57(4):592-613.
62. Giles MF, Rothwell PM. Systematic review and pooled analysis of published and unpublished validations of the ABCD and ABCD2 transient ischemic attack risk scores. *Stroke* 2010;41(4):667-73.

Table 1: Description of included reviews and prediction models

Systematic review	Giles 2010 ⁶²	Thompson 2014 ²¹	Siregar 2012 ²²	Damen ²⁹	Marques 2015 ²³	Tohira 2012 ²⁴	Klein 2013 ²⁵	Chalmers 2011 ²⁶	Ford 2010 ²⁷	Nassar 2014 ²⁸
Model	ABCD2 ³⁰	ESRS ³¹	EuroSCORE ³²	Framingham ³³	FRAX ³⁴	ISS ³⁵	MELD ³⁶	PSI ³⁷	RCRI ³⁸	SAPS 3 ³⁹
Population	Patients with TIA	Adults with a previous CVD event	Adult patients who underwent cardiac surgery under cardiopulmonary bypass	Men without previous CHD event	General population	Injured patients	Patients with liver cirrhosis but without hepatocellular carcinoma who underwent elective transjugular intrahepatic portosystemic shunts	Inpatients with community-acquired pneumonia	Patients aged >=50 years who underwent nonemergent noncardiac procedures	ICU patients
Geographical location (continent)	United States and UK (Europe and North America)	Canada, United States, Europe (Europe and North America)	Europe (Europe)	United States (North America)	Europe, Canada, Japan, United States, Australia (Europe, North America, Asia, Australia)	United States (North America)	United States (North America)	United States (North America)	United States (North America)	Worldwide (all continents)
Patient recruitment	1981-1998	1992-1995	1995	1971-1974	1980-1999	1968-1969	1991-1995	1989	1989-1994	2002
Predicted outcome	Stroke	Recurrent ischemic stroke, MI and vascular death	Mortality	CHD	Osteoporotic fractures	All-cause mortality	All-cause mortality	30-day hospital mortality	Major cardiac complications	Hospital mortality
Prediction horizon	2 days	1 year	30 days	10 years	10 years	3 months	3 months	30 days	1 year	90 days
Performance development study										
C-statistic	0.66 [95% CI 0.60, 0.71]	NR	0.7875	0.74	0.63	NR	NR	0.84	0.759 [SE 0.032]	0.848
OE ratio	NR*	NR*	NR*	NR*	NR*	NR*	NR*	NR*	NR*	1.00 [95% CI 0.98, 1.02]
Pooled performance validation studies										
Number of external validations included in analyses	16	11	22	23	30	34	14	24	23	27
C-statistic [95% CI]	0.66	0.60	0.79	0.68	0.66	0.86	0.64	0.80	0.69	0.83

	[0.61, 0.71]	[0.58, 0.62]	[0.77, 0.81]	[0.65, 0.71]	[0.63, 0.68]	[0.83, 0.88]	[0.59, 0.68]	[0.77, 0.82]	[0.65, 0.72]	[0.80, 0.85]
95% PI	[0.54, 0.77]	[0.57, 0.63]	[0.74, 0.83]	[0.56, 0.78]	[0.54, 0.76]	[0.62, 0.96]	[0.48, 0.77]	[0.64, 0.89]	[0.53, 0.81]	[0.66, 0.92]
OE ratio [95% CI]	NA	NA	0.54 [0.42, 0.68]	0.58 [0.45, 0.76]	1.10 [0.83, 1.47]	NA	NA	0.94 [0.83, 1.06]	2.70 [1.72, 4.25]	0.89 [0.77, 1.03]
95% PI	NA	NA	[0.19, 1.51]	[0.20, 1.74]	[0.31, 3.93]	NA	NA	[0.55, 1.60]	[0.35, 20.75]	[0.42, 1.91]

TIA: transient ischaemic attack, CVD: cardiovascular disease, CHD: coronary heart disease, ICU: intensive care unit, UK: United Kingdom, MI: myocardial infarction, NR: not reported, CI: confidence interval, PI: prediction interval, NA: not assessed.

*As the models are optimally fit in the development dataset, all OE ratios should be close to 1.

For peer review only

Figure legends

Figure 1: Flow chart of study selection.
SR: systematic review, IPD: individual participant data, MA: meta-analysis, NR: not reported, c: concordance, OE: observed expected.

Figure 2: Associations between study characteristics and logit c-statistic with regard to a reference category across 221 external validation studies and 10 different prediction models. Confidence intervals not including 0 are marked with an *. Figure S1 shows these differences on the original scale if we assume a c-statistic of 0.70 in the reference category. For example, for comparability of eligibility criteria, if we assume a c-statistic of 0.70 in the reference category (narrower), this would result in c-statistics of 0.74 [0.72, 0.77], 0.73 [0.66, 0.79], 0.77 [0.68, 0.84], and 0.77 [0.59,0.89] in the categories comparable, mixture, broader, and unclear, respectively.

Figure 3: C-statistic for categories of study design, pooled using univariable meta-regression analyses within each systematic review. N represents the number of external validation studies in a specific category. C diff represents the difference in c-statistic with regard to a reference category (indicated with 'ref').

Figure 4: Associations between study characteristics and ln OE ratio with regard to a reference category across 124 external validation studies and 6 different prediction models. Confidence intervals not including 0 are marked with an *. Figure S4 shows these differences on the original scale if we assume an OE ratio of 1.00 in the reference category. For example, for comparability of eligibility criteria, if we assume an OE ratio of 1.00 in the reference category (narrower), this would result in OE ratios of 0.83 [0.66, 1.05], 1.11 [0.72, 1.70], and 0.92 [0.54, 1.58] in the categories comparable, mixture, and broader, respectively.

Figure 5: Total OE ratio for categories of study design, pooled using univariable meta-regression analyses within each systematic review. N represents the number of external validation studies in a specific category. OE diff represents the difference in OE ratio with regard to a reference category (indicated with 'ref').

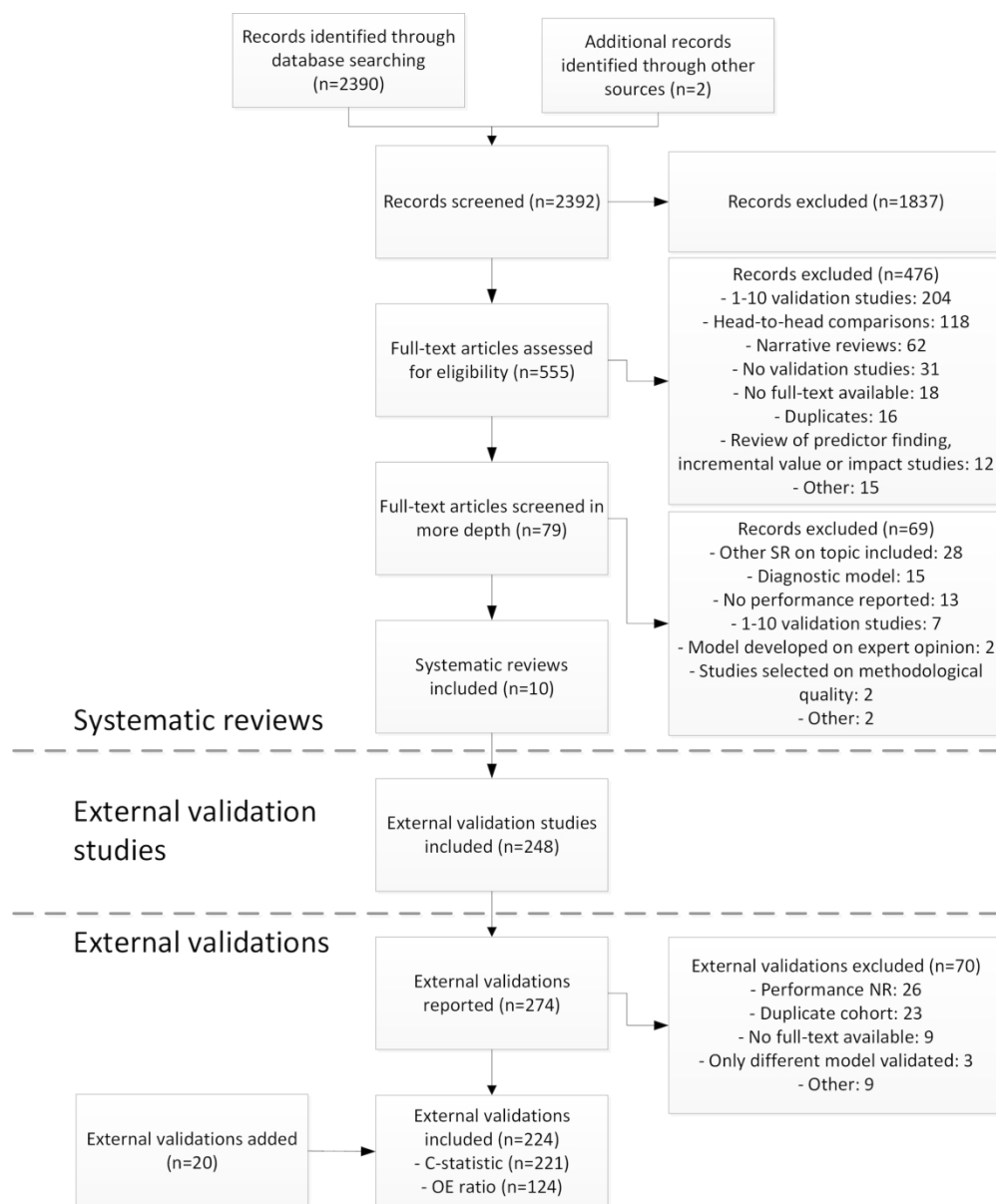


Figure 1: Flow chart of study selection. SR: systematic review, IPD: individual participant data, MA: meta-analysis, NR: not reported, c: concordance, OE: observed expected.

184x222mm (300 x 300 DPI)

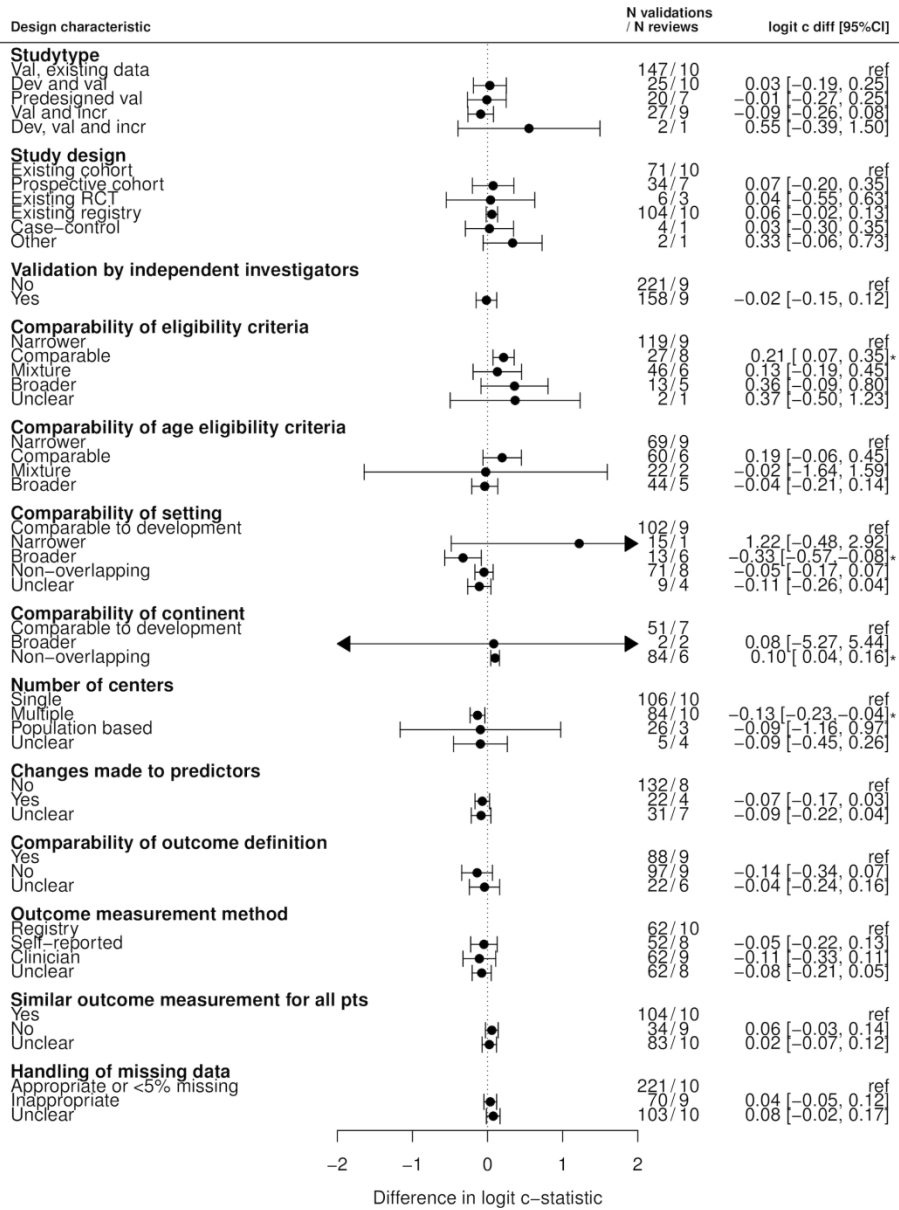


Figure 2: Associations between study characteristics and logit c-statistic with regard to a reference category across 221 external validation studies and 10 different prediction models. Confidence intervals not including 0 are marked with an *. Figure S1 shows these differences on the original scale if we assume a c-statistic of 0.70 in the reference category. For example, for comparability of eligibility criteria, if we assume a c-statistic of 0.70 in the reference category (narrower), this would result in c-statistics of 0.74 [0.72, 0.77], 0.73 [0.66, 0.79], 0.77 [0.68, 0.84], and 0.77 [0.59, 0.89] in the categories comparable, mixture, broader, and unclear, respectively.

167x222mm (300 x 300 DPI)

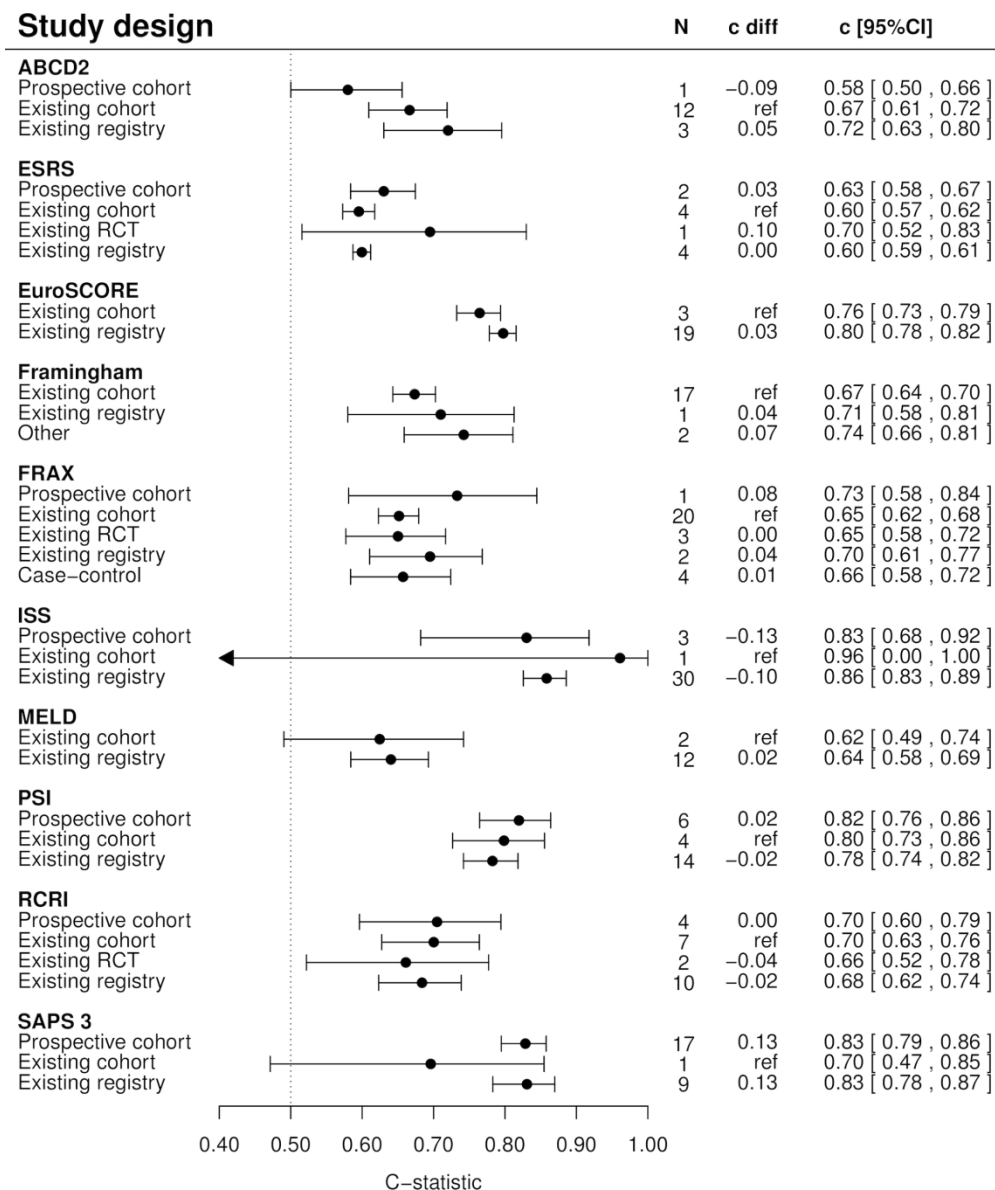


Figure 3: C-statistic for categories of study design, pooled using univariable meta-regression analyses within each systematic review. N represents the number of external validation studies in a specific category. C diff represents the difference in c-statistic with regard to a reference category (indicated with 'ref').

188x222mm (300 x 300 DPI)

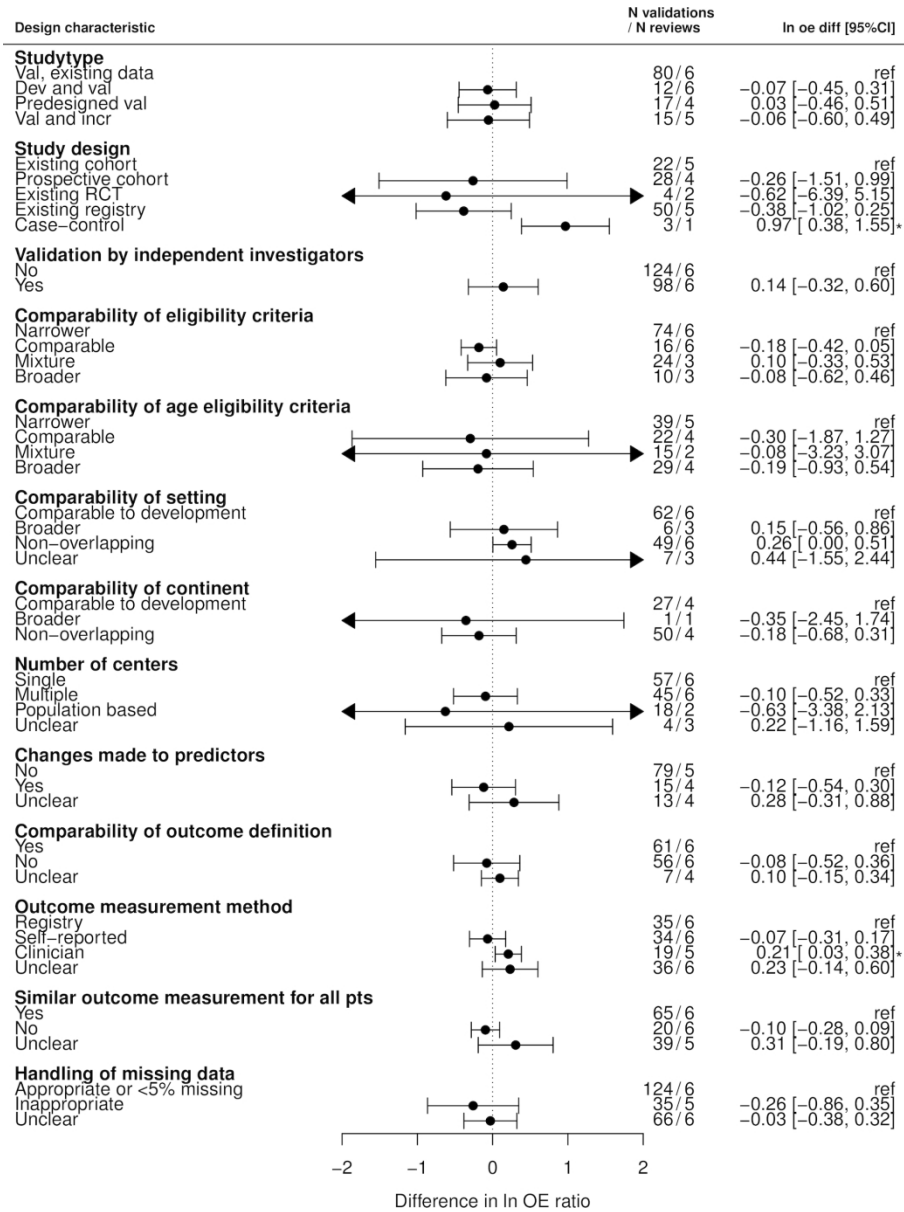


Figure 4: Associations between study characteristics and In OE ratio with regard to a reference category across 124 external validation studies and 6 different prediction models. Confidence intervals not including 0 are marked with an *. Figure S4 shows these differences on the original scale if we assume an OE ratio of 1.00 in the reference category. For example, for comparability of eligibility criteria, if we assume an OE ratio of 1.00 in the reference category (narrower), this would result in OE ratios of 0.83 [0.66, 1.05], 1.11 [0.72, 1.70], and 0.92 [0.54, 1.58] in the categories comparable, mixture, and broader, respectively.

167x222mm (300 x 300 DPI)

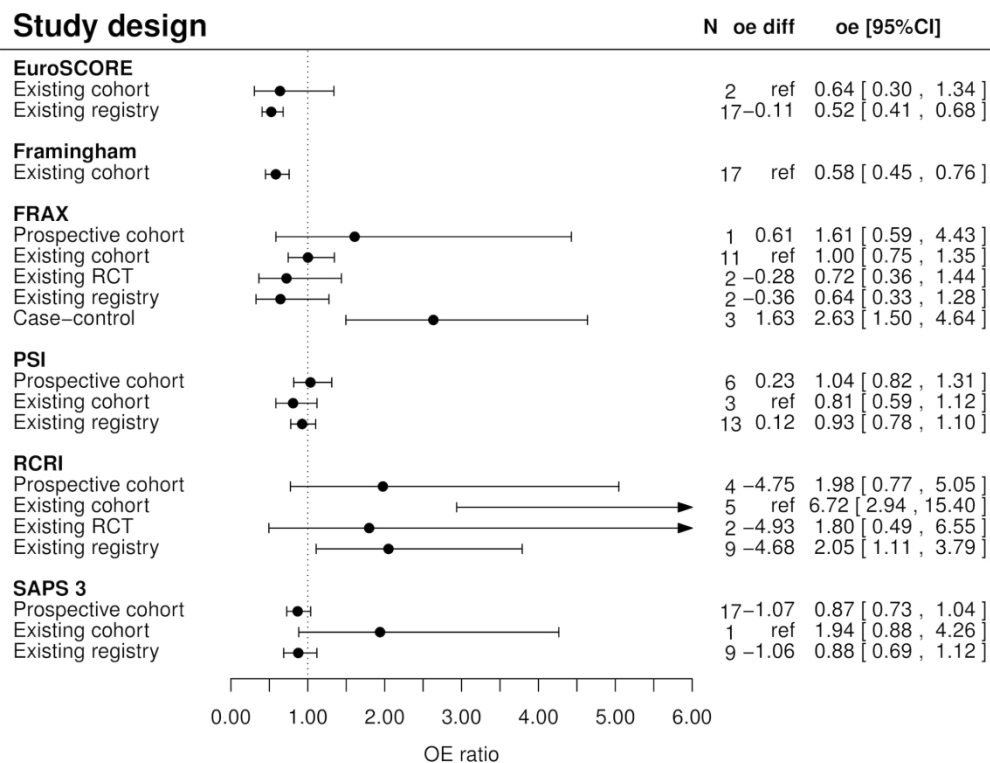


Figure 5: OE ratio for categories of study design, pooled using univariable meta-regression analyses within each systematic review. N represents the number of external validation studies in a specific category. OE diff represents the difference in OE ratio with regard to a reference category (indicated with 'ref').

190x142mm (300 x 300 DPI)

Damen et al. Empirical evidence on the impact of study characteristics and the performance of prediction models: a meta-epidemiological study

Supplement to “Empirical evidence on the impact of study characteristics and the performance of prediction models: a meta-epidemiological study”

Johanna A A G Damen, Thomas P A Debray, Romin Pajouheshnia, Johannes B Reitsma, Rob J P M Scholten, Karel G M Moons, Lotty Hooft

Content

Methods

- Supplement 1: Search string and selection criteria
- Supplement 2: Description of items extracted from studies and included in analyses
- Supplement 3: Statistical analyses

Tables and figures

- Table S1: Description of study characteristics and quality of reporting within each systematic review
- Figure S1: Associations between categorical variables and c-statistic
- Figure S2: Associations between continuous variables and c-statistic
- Figure S3: C-statistic in categories of study characteristics within each systematic review
- Figure S4: Associations between categorical variables and total OE ratio
- Figure S5: Associations between continuous variables and total OE ratio
- Figure S6: Total OE ratio in categories of study characteristics within each systematic review

Damen et al. Empirical evidence on the impact of study characteristics and the performance of prediction models: a meta-epidemiological study

Supplement 1: Search string and selection criteria

The systematic reviews included in this article have been selected from a database with reviews of prediction models. To construct this database, the following search string was used to search Pubmed (last update: October 16th, 2018):

(clinical prediction[ti] OR
risk calculator*[ti] OR
risk index[ti] OR
risk indices[ti] OR
risk model*[ti] OR
risk prediction[ti] OR
risk score*[ti] OR
risk stratification[ti] OR
predictive model*[ti] OR
prediction model*[ti] OR
prediction rule*[tiab] OR
prognostic index[ti] OR
prognostic indices[ti] OR
prognostic model*[ti] OR
scoring system*[ti]) AND
(review[Publication Type] OR
review[ti] OR
critical appraisal[ti] OR
Bibliography[Publication Type] OR
Meta-analysis[Publication Type]) NOT
(Editorial[Publication Type] OR
Letter[Publication Type] OR
News[Publication Type])

All references identified with this search string were screened for eligibility to the database by one reviewer (Gary Collins) based on title, abstract and, if necessary, full text. The following in- and exclusion criteria were applied:

Inclusion criteria

- Review (narrative or systematic) of prediction models
- Both diagnostic and prognostic prediction models
- Validation studies in which multiple prediction models have been validated (head-to-head comparison)
- Review of impact studies of prediction models
- Protocol of systematic review of prediction models

Exclusion criteria

- Comments, editorials

Supplement 2: Description of items extracted from studies and included in analyses

Item	Extracted from studies	Categorization / handling in analyses	Description / examples
Validated model	ABCD2, ESRS, EuroSCORE, Framingham Wilson, FRAX, ISS, MELD, PSI, RCRI, SAPS 3	-	-
Study type	Predesigned validation study	Predesigned validation study	Study designed with the aim of validating a prediction model
	Validation study using existing data	Validation study using existing data	Study in which a prediction model is validated using a dataset collected for a different purpose than validating the model
	Development of new model and validation of different model	Development of new model and validation of different model	Study in which a model is developed and a model is validated
	Validation and incremental value	Validation and incremental value	Study in which a model is validated and in which the added value of one or more predictors is assessed
	Development, validation, and incremental value study	Development, validation, and incremental value study	Combination of the two above
Independent investigators	Yes	Yes	None of the authors of the development study was listed as author in the external validation study
	No	No	One or more of the authors of the development study was listed as author in the external validation study
Study design	Prospective cohort	Prospective cohort	
	Existing cohort	Existing cohort	
	Existing RCT	Existing RCT	
	Existing registry / medical records	Existing registry	
	Case-control	Case-control	
	Other (specify)	Other	
Eligibility criteria for participants	Copy/paste eligibility criteria of validation study	Comparable	Eligibility criteria comparable to development study

Damen et al. Empirical evidence on the impact of study characteristics and the performance of prediction models: a meta-epidemiological study

Item	Extracted from studies	Categorization / handling in analyses	Description / examples
		Narrower	People included in the development study excluded in the validation study
		Broader	People excluded in the development study included in the validation study
		Mixture	Combination of narrower and broader
		Unclear	
Setting	Primary care Secondary care Tertiary care Population based Screening Mixed Unclear	Comparable	Same setting as development study
		Broader	Same setting as development study, and participants from additional settings recruited
		Non-overlapping	Setting in development study differs from validation study
		Unclear	
Study dates	Start year of recruitment End year of recruitment	Continuous, standardized per systematic review	
Prediction horizon	Time period for which predictions were made, eg, 10 years.	Continuous, standardized per systematic review	
Geographical location	Country and continent	Comparable	Model validated in the same continent as the development study
		Broader	Model validated in the same and additional continents as the development study
		Non-overlapping	Model validated in a different continent than the development study
Number of centres	Number of centres (numerical)	Single	
		Multiple	
		Population based	Participants not recruited at medical centres, but, for example, from a specific geographic area (eg, all individuals living in Framingham, US)
		Unclear	
Case-mix: age mean	Mean and SD of age of	Continuous, standardized per	

Damen et al. Empirical evidence on the impact of study characteristics and the performance of prediction models: a meta-epidemiological study

Item	Extracted from studies	Categorization / handling in analyses	Description / examples
and sd	participants included in the study, or other available information about age distribution	systematic review	
Case-mix: gender	Percentage of men included in a study	Continuous, standardized per systematic review	
Predictors	Were predictors deleted from the model, or were predictors substituted with different predictors.	Yes	Changes made to predictors
		No	No changes made to predictors
		Unclear	
Predicted outcome	Full definition, including ICD-codes	Comparable	Outcome definition comparable to development study
		Not comparable	Outcome definition not comparable to development study
		Unclear	
Outcome - measurement method	Measurement method (eg, self-reported, interviews, expert panel), differences in outcome measurement between participants in the study	Yes	Outcome measurement similar for all participant
		No	Systematic differences in outcome measurement between participants
		Unclear	
Missing data	Number of participants with missing data, method of handling missing data	Appropriate	Missing data handled using multiple imputation, or <5% missing data (arbitrary cut-off)
		Inappropriate	Missing data not handled using multiple imputation (eg, complete-case analysis, mean imputation), and >=5% missing data
		Unclear	Unclear handling of missing data, and >=5% missing data
Number of participants		Continuous, standardized per systematic review	
Number of events		Continuous, standardized per systematic review	

Damen et al. Empirical evidence on the impact of study characteristics and the performance of prediction models: a meta-epidemiological study

Item	Extracted from studies	Categorization / handling in analyses	Description / examples
Model updating	Was the model altered before validating it, eg, using intercept recalibration.	NA	
Performance - c-statistic	C-statistic, AUC, 95% confidence intervals or SE	Logit transformation ¹	
Performance - total OE ratio	OE ratio, predicted risks, presence of calibration plots or tables, 95% confidence intervals or SE	Ln transformation ¹	

SD: standard deviation, NA: not applicable, C-statistic: concordance statistic, AUC: area under the receiver operating curve, SE: standard error, OE ratio: observed expected ratio.

Information regarding c-statistics and total OE ratios when not reported was sometimes restored from other information reported in the paper. If the precision of the c-statistic was not reported, we estimated this from the c-statistic and sample size of the study, using the formula described by Newcombe and Hanley.^{2,3} Various equations were used to estimate the standard error of the OE ratio, depending on which information was reported. All equations (as numbered) are described in the appendix of Debray et al.⁴ If the SE of the OE ratio was reported, we used equation 16 to estimate the SE of $\ln(\text{OE})$, if the observed event risk (P_o), the expected event risk (P_e), and the SE of P_o were reported, we used equation 51, and if only P_o and P_e were reported we used equation 27.

Supplement 3: Statistical analyses

First we pooled the total OE ratio and c-statistic within each systematic review. Based on previous recommendations,¹⁴ we pooled the log OE ratio and logit c-statistic using random-effects meta-analysis accounting for the presence of between-study heterogeneity, weighted by the inverse of the variance. We calculated 95% confidence intervals (CI) and (approximate) 95% prediction intervals (PI) to quantify uncertainty and the presence of between-study heterogeneity. The Hartung-Knapp-Sidik-Jonkman (HKSJ) method was used when calculating 95% CIs.⁵ The 95% PI was calculated using the equation described previously.⁴ The CI indicates the precision of the summary performance estimate and the PI provides boundaries on the likely performance in future model validation studies that are comparable to the studies included in the meta-analysis, and can thus be seen as an indication of model generalizability.⁶

To study the possible association between study characteristics and predictive performance, we used a two-stepped approach. In the first stage, we fitted a univariable meta-regression model (ie, a separate model for every study characteristic) within every systematic review, with the logit c-statistic or log OE ratio as outcome variable. This model was fitted with intercept term. Therefore, the effect estimates obtained from this meta-regression model indicate the difference in logit c-statistic or log OE ratio between a certain category of a study characteristic and a chosen reference category of that characteristic. As a reference category, we chose the category that was present in the highest number of systematic reviews allowing the inclusion of as many data as possible.

In the second stage, these effect estimates were pooled with a random effects meta-analysis model. This reflected the influence of the study characteristic on model performance over all systematic reviews. For continuous study characteristics, the intercept term and beta-coefficient from the first stage were jointly pooled across reviews using bivariate meta-analysis.^{4 6} For categorical study characteristics the data available were not sufficient for the complexity of a multivariate model, so every category was pooled in a separate (univariate) meta-analysis.

As the estimates obtained with this approach are on the transformed scale (ie, the difference in logit c-statistic or log OE ratio between one category and the reference category), we transformed these back assuming a c-statistic of 0.7 or an OE ratio of 1.00 in the reference category. Also, we performed a second analysis where we again fitted a univariable meta-regression model, with the logit c-statistic or log OE ratio as outcome variable, but now without intercept term. This analysis enables the calculation

Damen et al. Empirical evidence on the impact of study characteristics and the performance of prediction models: a meta-epidemiological study

of an effect estimate for every category of a study characteristic and to back transform this to the original scale, yielding a pooled c-statistic or pooled OE ratio for each category of a study characteristic.

We planned to perform multivariable analyses to assess the association between various study characteristics in combination and the performance of prediction models, but due to the paucity of data we were not able to do so.

All analyses were performed in R version 3.3.2,⁷ using the packages metafor,⁸ mvmeta,⁹ metamisc,¹⁰ and lme4.¹¹

For peer review only

Damen et al. Empirical evidence on the impact of study characteristics and the performance of prediction models: a meta-epidemiological study
Table S1: Description of study characteristics and quality of reporting within each systematic review
Categorical variables

	ABCD2	ESRS	EuroSCORE	Framingham	FRAX	ISS	MELD	PSI	RCRI	SAPS 3
Studytype										
Validation study using existing data	9 (56%)	7 (64%)	21 (95%)	18 (78%)	26 (87%)	24 (71%)	10 (71%)	16 (67%)	11 (48%)	8 (30%)
Development of new model and validation of different model	2 (12%)	2 (18%)	1 (5%)	2 (9%)	1 (3%)	4 (12%)	3 (21%)	2 (8%)	2 (9%)	6 (22%)
Development, validation, and incremental value study	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	2 (6%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)
Predesigned validation study	1 (6%)	1 (9%)	0 (0%)	0 (0%)	1 (3%)	1 (3%)	0 (0%)	5 (21%)	1 (4%)	10 (37%)
Validation and incremental value	4 (25%)	1 (9%)	0 (0%)	3 (13%)	2 (7%)	3 (9%)	1 (7%)	1 (4%)	9 (39%)	3 (11%)
Study design										
Existing cohort	12 (75%)	4 (36%)	3 (14%)	20 (87%)	20 (67%)	1 (3%)	2 (14%)	4 (17%)	7 (30%)	1 (4%)
Prospective cohort	1 (6%)	2 (18%)	0 (0%)	0 (0%)	1 (3%)	3 (9%)	0 (0%)	6 (25%)	4 (17%)	17 (63%)
Existing RCT	0 (0%)	1 (9%)	0 (0%)	0 (0%)	3 (10%)	0 (0%)	0 (0%)	0 (0%)	2 (9%)	0 (0%)
Existing registry	3 (19%)	4 (36%)	19 (86%)	1 (4%)	2 (7%)	30 (88%)	12 (86%)	14 (58%)	10 (43%)	9 (33%)
Case-control	0 (0%)	0 (0%)	0 (0%)	0 (0%)	4 (13%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)
Other	0 (0%)	0 (0%)	0 (0%)	2 (9%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)
Validation by independent investigators										
No	6 (38%)	3 (27%)	3 (14%)	10 (43%)	17 (57%)	2 (6%)	0 (0%)	7 (29%)	1 (4%)	2 (7%)
Yes	10 (62%)	8 (73%)	19 (86%)	13 (57%)	13 (43%)	32 (94%)	14 (100%)	17 (71%)	22 (96%)	25 (93%)
Comparability of eligibility criteria										
Narrower	6 (38%)	2 (18%)	18 (82%)	18 (78%)	28 (93%)	22 (65%)	0 (0%)	4 (17%)	4 (17%)	20 (74%)
Comparable	4 (25%)	0 (0%)	2 (9%)	3 (13%)	2 (7%)	5 (15%)	0 (0%)	1 (4%)	3 (13%)	7 (26%)
Mixture	5 (31%)	9 (82%)	0 (0%)	2 (9%)	0 (0%)	3 (9%)	0 (0%)	16 (67%)	11 (48%)	0 (0%)
Broader	1 (6%)	0 (0%)	2 (9%)	0 (0%)	0 (0%)	2 (6%)	0 (0%)	3 (12%)	5 (22%)	0 (0%)
Non-overlapping	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	14 (100%)	0 (0%)	0 (0%)	0 (0%)
Unclear	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	2 (6%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)
Comparability of age eligibility criteria										
Narrower	1 (6%)	2 (18%)	0 (0%)	17 (74%)	9 (30%)	13 (38%)	10 (71%)	2 (8%)	1 (4%)	17 (63%)
Comparable	15 (94%)	0 (0%)	22 (100%)	0 (0%)	1 (3%)	21 (62%)	4 (29%)	15 (62%)	4 (17%)	4 (15%)
Mixture	0 (0%)	0 (0%)	0 (0%)	6 (26%)	16 (53%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)
Broader	0 (0%)	9 (82%)	0 (0%)	0 (0%)	4 (13%)	0 (0%)	0 (0%)	7 (29%)	18 (78%)	6 (22%)
Setting										

Damen et al. Empirical evidence on the impact of study characteristics and the performance of prediction models: a meta-epidemiological study

	ABCD2	ESRS	EuroSCORE	Framingham	FRAX	ISS	MELD	PSI	RCRI	SAPS 3
Primary care	3 (19%)	0 (0%)	0 (0%)	7 (30%)	3 (10%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)
Secondary care	12 (75%)	9 (82%)	16 (84%)	0 (0%)	5 (17%)	18 (82%)	6 (75%)	18 (90%)	12 (86%)	4 (40%)
Tertiary care	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)
Population based	0 (0%)	1 (9%)	0 (0%)	15 (65%)	17 (59%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)
Screening	0 (0%)	0 (0%)	0 (0%)	1 (4%)	1 (3%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)
Mixed	1 (6%)	1 (9%)	2 (11%)	0 (0%)	0 (0%)	4 (18%)	1 (12%)	2 (10%)	2 (14%)	2 (20%)
Unclear	0 (0%)	0 (0%)	1 (5%)	0 (0%)	3 (10%)	0 (0%)	1 (12%)	0 (0%)	0 (0%)	4 (40%)
Comparability of setting										
Comparable	1 (6%)	0 (0%)	16 (73%)	15 (65%)	17 (57%)	18 (53%)	6 (43%)	18 (75%)	9 (39%)	4 (15%)
Narrower	15 (94%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)
Broader	0 (0%)	0 (0%)	2 (9%)	0 (0%)	0 (0%)	4 (12%)	1 (7%)	2 (8%)	2 (9%)	2 (7%)
Non-overlapping	0 (0%)	11 (100%)	3 (14%)	8 (35%)	10 (33%)	12 (35%)	6 (43%)	4 (17%)	12 (52%)	17 (63%)
Unclear	0 (0%)	0 (0%)	1 (5%)	0 (0%)	3 (10%)	0 (0%)	1 (7%)	0 (0%)	0 (0%)	4 (15%)
Continent										
Africa	0 (0%)	0 (0%)	1 (5%)	0 (0%)	0 (0%)	1 (3%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)
Asia	0 (0%)	2 (18%)	4 (18%)	2 (9%)	4 (13%)	7 (21%)	1 (7%)	3 (12%)	2 (9%)	5 (19%)
Australia	0 (0%)	0 (0%)	1 (5%)	0 (0%)	5 (17%)	2 (6%)	0 (0%)	3 (12%)	1 (4%)	1 (4%)
Europe	8 (50%)	7 (64%)	10 (45%)	10 (43%)	9 (30%)	7 (21%)	7 (50%)	11 (46%)	13 (57%)	9 (33%)
North America	8 (50%)	1 (9%)	5 (23%)	11 (48%)	11 (37%)	17 (50%)	3 (21%)	6 (25%)	6 (26%)	3 (11%)
South America	0 (0%)	0 (0%)	1 (5%)	0 (0%)	0 (0%)	0 (0%)	3 (21%)	0 (0%)	0 (0%)	9 (33%)
Combination	0 (0%)	1 (9%)	0 (0%)	0 (0%)	1 (3%)	0 (0%)	0 (0%)	1 (4%)	1 (4%)	0 (0%)
Comparability of continent										
Comparable	0 (0%)	0 (0%)	10 (45%)	11 (48%)	0 (0%)	17 (50%)	3 (21%)	6 (25%)	6 (26%)	0 (0%)
Narrower	16 (100%)	8 (73%)	0 (0%)	0 (0%)	30 (100%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	27 (100%)
Broader	0 (0%)	1 (9%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	1 (4%)	1 (4%)	0 (0%)
Non-overlapping	0 (0%)	2 (18%)	12 (55%)	12 (52%)	0 (0%)	17 (50%)	11 (79%)	17 (71%)	16 (70%)	0 (0%)
Number of centres										
Single	9 (56%)	4 (36%)	12 (55%)	3 (13%)	5 (17%)	18 (53%)	12 (86%)	14 (58%)	17 (74%)	13 (48%)
Multiple	6 (38%)	7 (64%)	9 (41%)	6 (26%)	9 (30%)	15 (44%)	2 (14%)	10 (42%)	6 (26%)	14 (52%)
Population based	1 (6%)	0 (0%)	0 (0%)	12 (52%)	15 (50%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)
Unclear	0 (0%)	0 (0%)	1 (5%)	2 (9%)	1 (3%)	1 (3%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)

Damen et al. Empirical evidence on the impact of study characteristics and the performance of prediction models: a meta-epidemiological study

	ABCD2	ESRS	EuroSCORE	Framingham	FRAX	ISS	MELD	PSI	RCRI	SAPS 3
Changes made to predictors										
No	16 (100%)	10 (91%)	13 (59%)	23 (100%)	20 (67%)	21 (62%)	12 (86%)	13 (54%)	17 (74%)	26 (96%)
Yes	0 (0%)	0 (0%)	5 (23%)	0 (0%)	10 (33%)	0 (0%)	0 (0%)	5 (21%)	2 (9%)	0 (0%)
Unclear	0 (0%)	1 (9%)	4 (18%)	0 (0%)	0 (0%)	13 (38%)	2 (14%)	6 (25%)	4 (17%)	1 (4%)
Comparability of outcome definition										
No	8 (50%)	3 (27%)	9 (41%)	13 (57%)	16 (53%)	7 (21%)	14 (100%)	5 (21%)	5 (22%)	24 (89%)
Yes	3 (19%)	8 (73%)	11 (50%)	4 (17%)	13 (43%)	19 (56%)	0 (0%)	18 (75%)	18 (78%)	3 (11%)
Unclear	5 (31%)	0 (0%)	2 (9%)	6 (26%)	1 (3%)	8 (24%)	0 (0%)	1 (4%)	0 (0%)	0 (0%)
Outcome measurement method										
Self-reported	3 (19%)	8 (73%)	2 (9%)	15 (65%)	18 (60%)	0 (0%)	1 (7%)	6 (25%)	1 (4%)	2 (7%)
Clinician	6 (38%)	1 (9%)	2 (9%)	2 (9%)	1 (3%)	4 (12%)	0 (0%)	2 (8%)	10 (43%)	6 (22%)
Registry	3 (19%)	2 (18%)	5 (23%)	4 (17%)	8 (27%)	13 (38%)	3 (21%)	9 (38%)	6 (26%)	9 (33%)
Unclear	4 (25%)	0 (0%)	13 (59%)	2 (9%)	3 (10%)	17 (50%)	10 (71%)	7 (29%)	6 (26%)	10 (37%)
Similar outcome measurement for all patients										
Yes	9 (56%)	3 (27%)	11 (50%)	6 (26%)	14 (47%)	18 (53%)	2 (14%)	12 (50%)	13 (57%)	16 (59%)
No	1 (6%)	5 (45%)	1 (5%)	16 (70%)	3 (10%)	3 (9%)	0 (0%)	3 (12%)	4 (17%)	1 (4%)
Unclear	6 (38%)	3 (27%)	10 (45%)	1 (4%)	13 (43%)	13 (38%)	12 (86%)	9 (38%)	6 (26%)	10 (37%)
Method for handling of missing data										
Complete case analysis	4 (25%)	8 (73%)	3 (14%)	11 (48%)	19 (63%)	18 (53%)	7 (50%)	1 (4%)	5 (22%)	8 (30%)
Mean/median imputation	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	1 (4%)
Multiple imputation	1 (6%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)
NA	3 (19%)	1 (9%)	1 (5%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	1 (4%)	1 (4%)	2 (7%)
Other	0 (0%)	0 (0%)	1 (5%)	0 (0%)	2 (7%)	0 (0%)	1 (7%)	8 (33%)	0 (0%)	7 (26%)
Unclear	8 (50%)	2 (18%)	17 (77%)	12 (52%)	9 (30%)	16 (47%)	6 (43%)	14 (58%)	17 (74%)	9 (33%)
Handling of missing data										
Appropriate or <5% missing	8 (50%)	4 (36%)	2 (9%)	3 (13%)	4 (13%)	6 (18%)	4 (29%)	4 (17%)	7 (30%)	6 (22%)
Inappropriate	1 (6%)	5 (45%)	3 (14%)	8 (35%)	17 (57%)	13 (38%)	4 (29%)	8 (33%)	0 (0%)	12 (44%)
Unclear	7 (44%)	2 (18%)	17 (77%)	12 (52%)	9 (30%)	15 (44%)	6 (43%)	12 (50%)	16 (70%)	9 (33%)

Damen et al. Empirical evidence on the impact of study characteristics and the performance of prediction models: a meta-epidemiological study

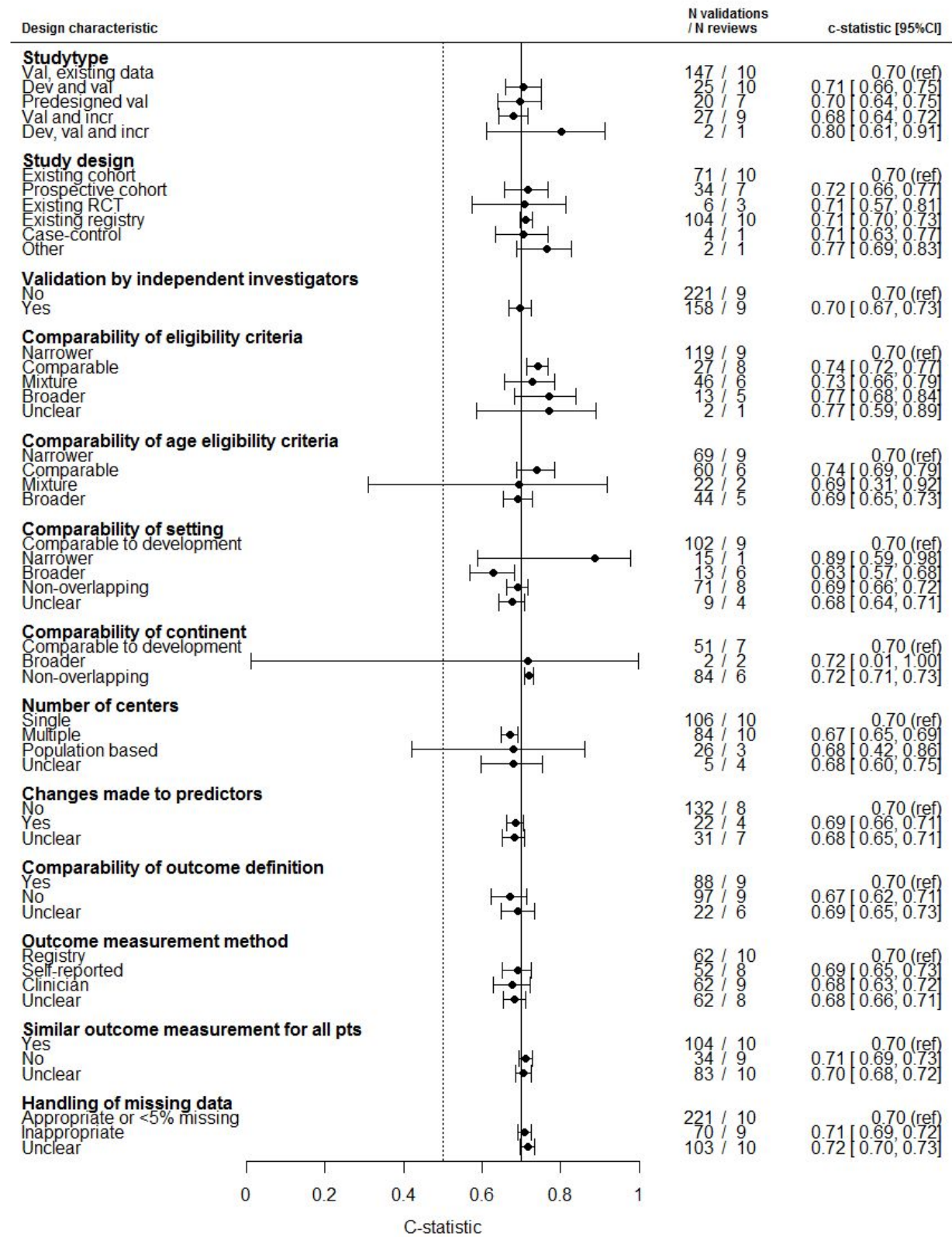
Continuous variables

	ABCD2	ESRS	EuroSCORE	Framingham	FRAX	ISS	MELD	PSI	RCRI	SAPSIII
Year start recruitment	2002 (2000-2003) NR=0	2007 (2004-2007) NR=0	1998 (1995-2001) NR=1	1989 (1983-1994) NR=0	1994 (1990-1998) NR=3	1996 (1993-1998) NR=1	2000 (1998-2004) NR=0	2000 (1998-2002) NR=0	2000 (1994-2002) NR=4	2006 (2006-2007) NR=0
Year end recruitment	2005 (2003-2007) NR=0	2008 (2006-2008) NR=0	2002 (1999-2005) NR=1	1993 (1988-1998) NR=0	1997 (1993-2006) NR=8	2000 (1996-2003) NR=2	2006 (2004-2007) NR=0	2002 (2000-2003) NR=0	2002 (2000-2005) NR=4	2007 (2006-2009) NR=0
Percentage missings	0.95 (0.00-5.00) NR=7	5.12 (1.99-17.80) NR=2	6.40 (1.50-11.83) NR=18	4.90 (2.70-9.80) NR=18	30.25 (2.75-33.80) NR=16	9.05 (2.40-14.65) NR=20	4.05 (2.73-10.93) NR=8	0.52 (0.07-9.26) NR=18	1.00 (0.09-1.91) NR=16	5.85 (0.52-18.93) NR=15
Number of participants	304 (204-691) NR=0	1257 (712-2594) NR=0	1730 (873-4518) NR=2	2399 (928-4609) NR=0	2210 (889-6586) NR=0	2590 (960-20713) NR=0	418 (118-483) NR=0	730 (326-970) NR=1	496 (180-1480) NR=0	864 (485-1856) NR=0
Number of events	9 (3-18) NR=0	92 (60-134) NR=0	36 (13-87) NR=2	92 (72-160) NR=1	250 (86-581) NR=0	256 (113-1660) NR=2	49 (22-112) NR=0	54 (28-111) NR=1	31 (14-76) NR=0	180 (124-311) NR=1
Age mean	67.4 (64.1-70.0) NR=5	68.3 (67.1-71.5) NR=3	63.9 (62.5-65.2) NR=2	54.6 (50.9-58.3) NR=2	66.8 (63.0-71.3) NR=1	38.1 (32.4-41.3) NR=10	51.8 (49.1-53.0) NR=0	66.2 (64.0-69.3) NR=2	67.8 (66.0-71.9) NR=2	62.2 (60.8-64.8) NR=1
Age sd	13.8 (13.0-14.9) NR=5	12.4 (12.0-13.0) NR=1	9.3 (9.0-10.6) NR=8	7.3 (4.1-9.4) NR=0	8.3 (5.9-9.8) NR=0	20.9 (18.1-24.8) NR=2	10.0 (9.6-12.0) NR=1	17.8 (17.0-20.1) NR=3	10.0 (8.8-12.5) NR=4	17.0 (15.4-19.0) NR=3
Gender percentage men	47 (45-53) NR=2	57 (55-59) NR=1	77 (71-79) NR=1	100 (100-100) NR=0	0 (0-0) NR=0	71 (64-75) NR=11	68 (63-69) NR=0	57 (53-64) NR=1	67 (52-76) NR=0	59 (55-64) NR=0

Values represent median (IQR), number of missing values.

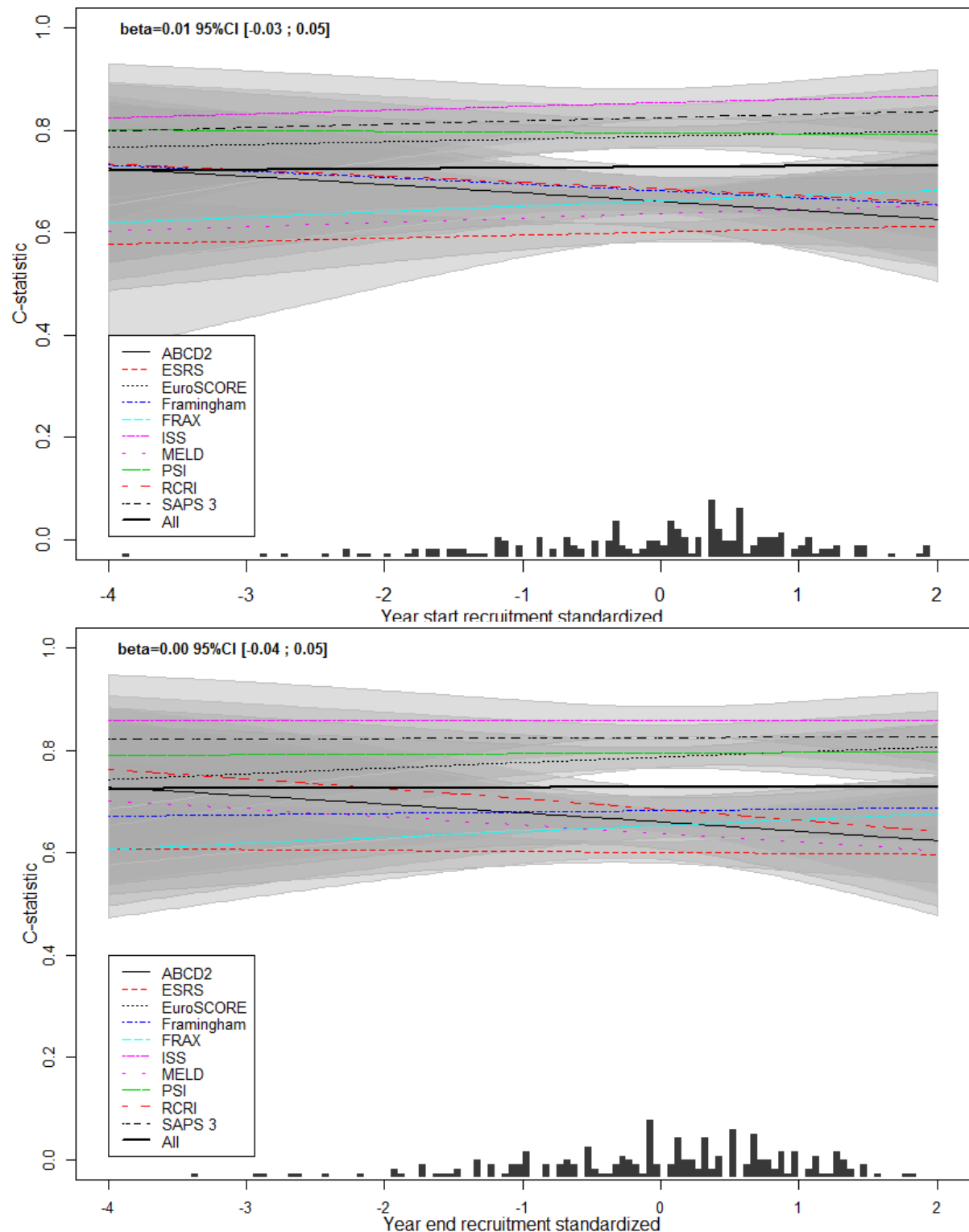
Damen et al. Empirical evidence on the impact of study characteristics and the performance of prediction models: a meta-epidemiological study

Figure S1: Associations between categorical variables and c-statistic

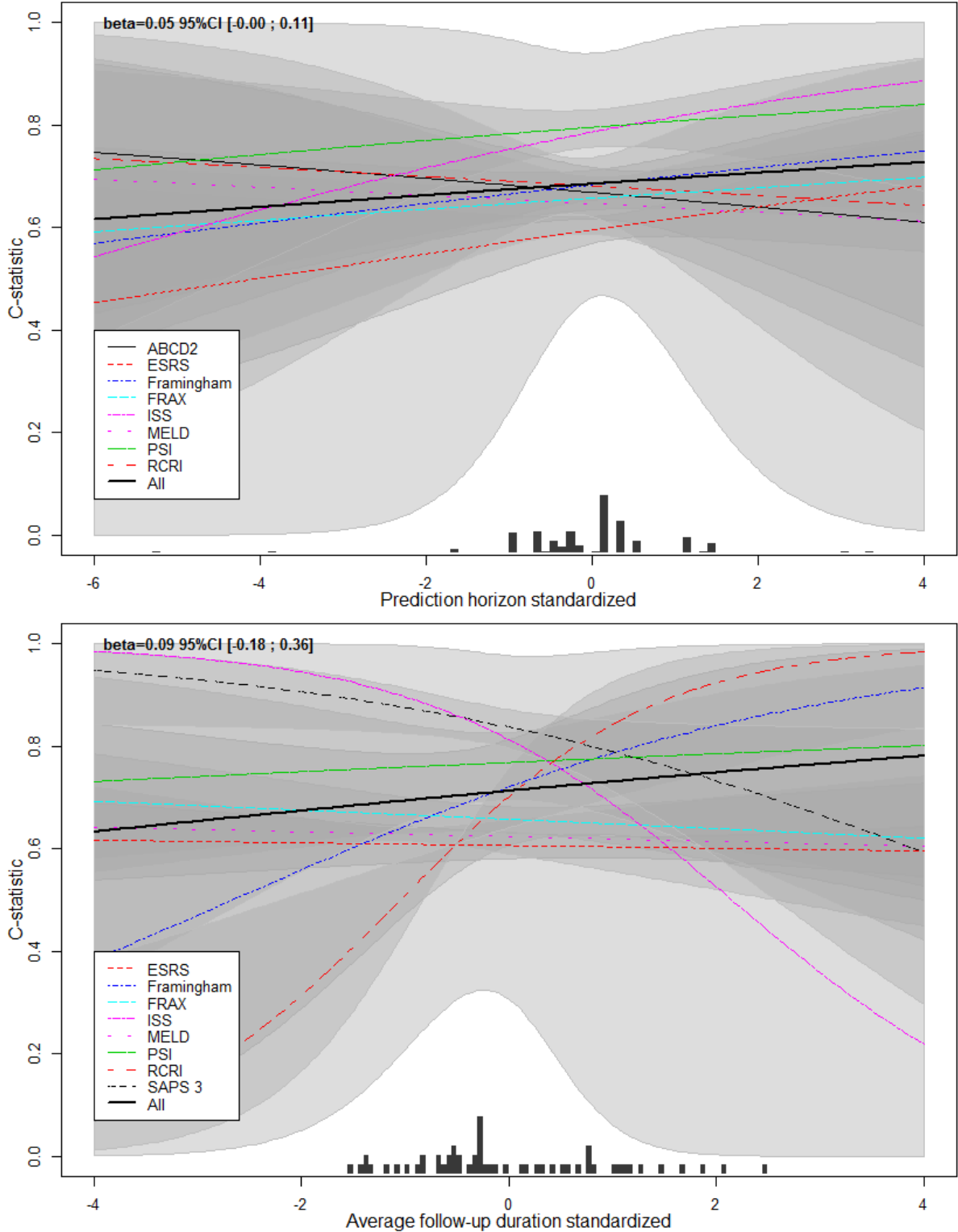


Damen et al. Empirical evidence on the impact of study characteristics and the performance of prediction models: a meta-epidemiological study

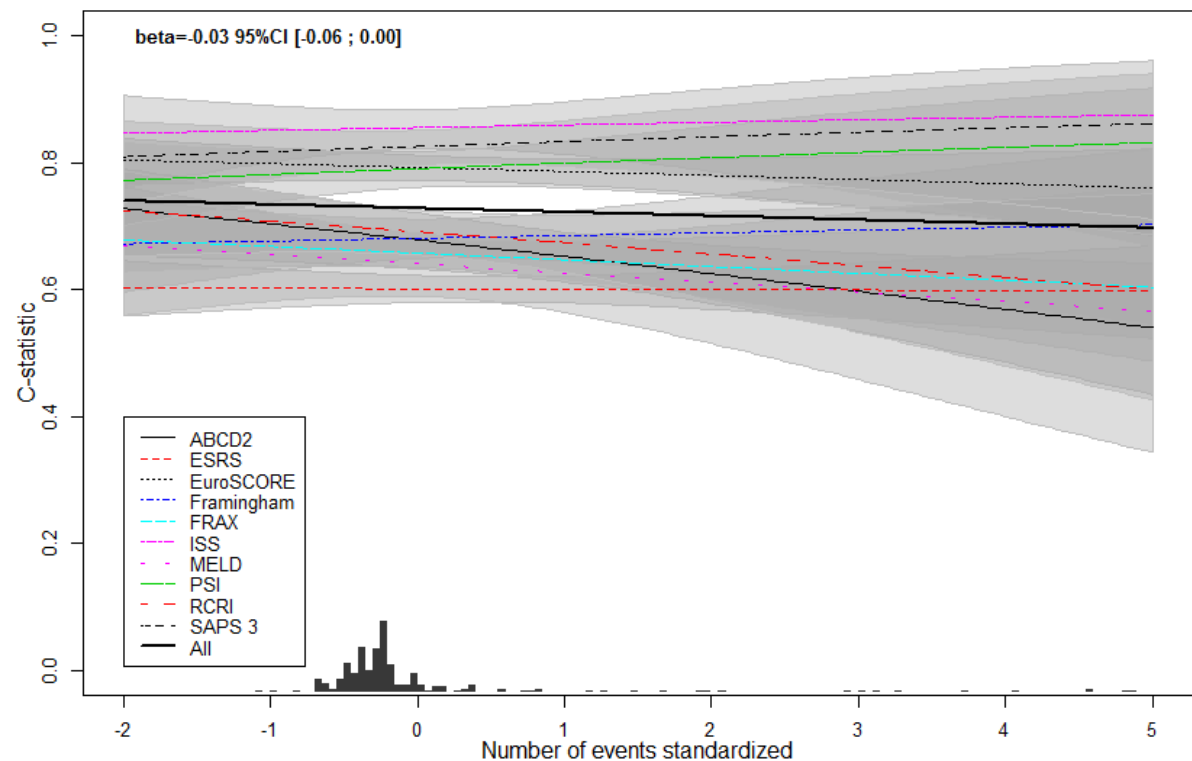
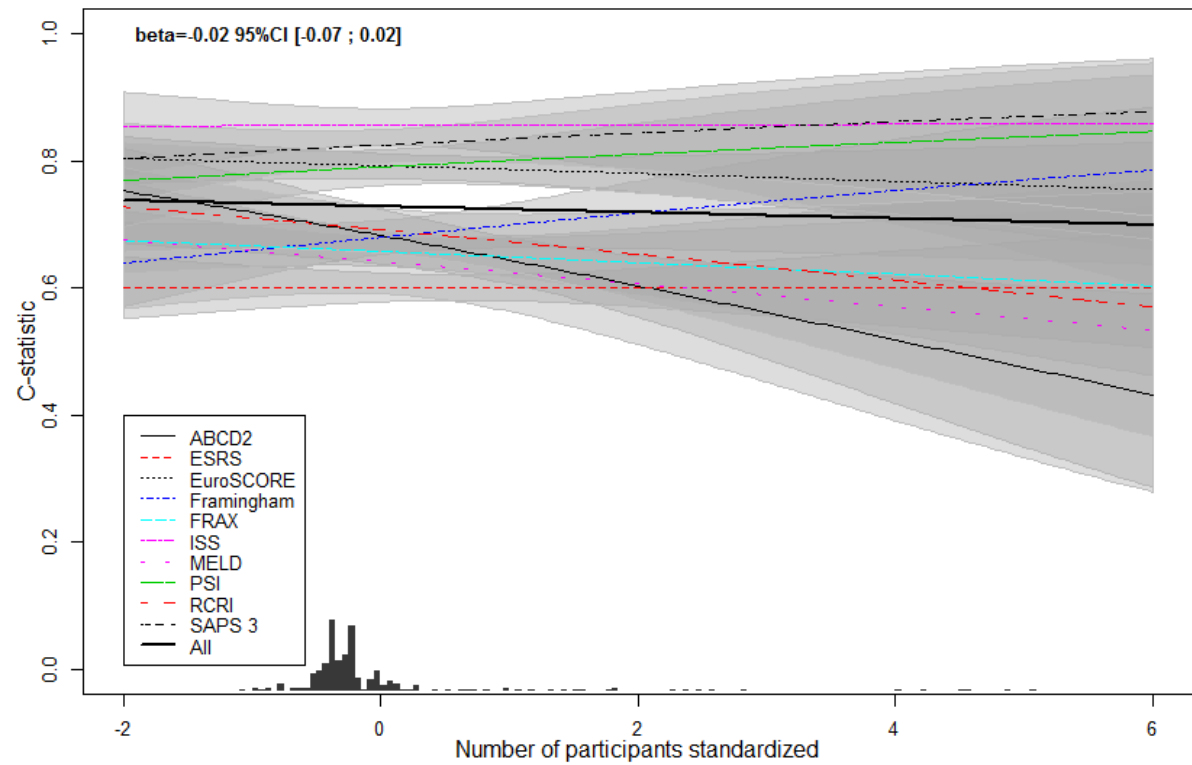
Figure S2: Associations between continuous variables and c-statistic



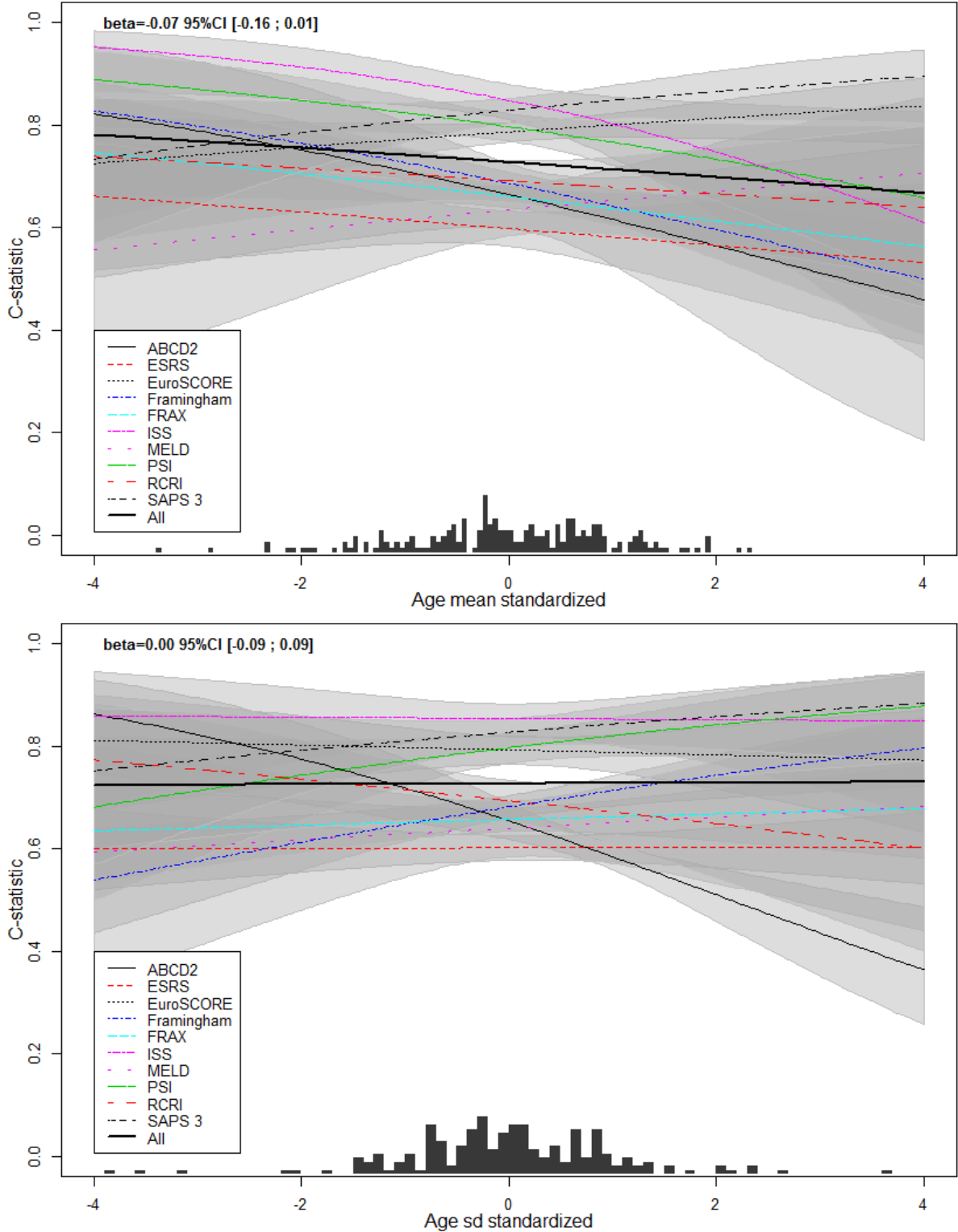
Damen et al. Empirical evidence on the impact of study characteristics and the performance of prediction models: a meta-epidemiological study



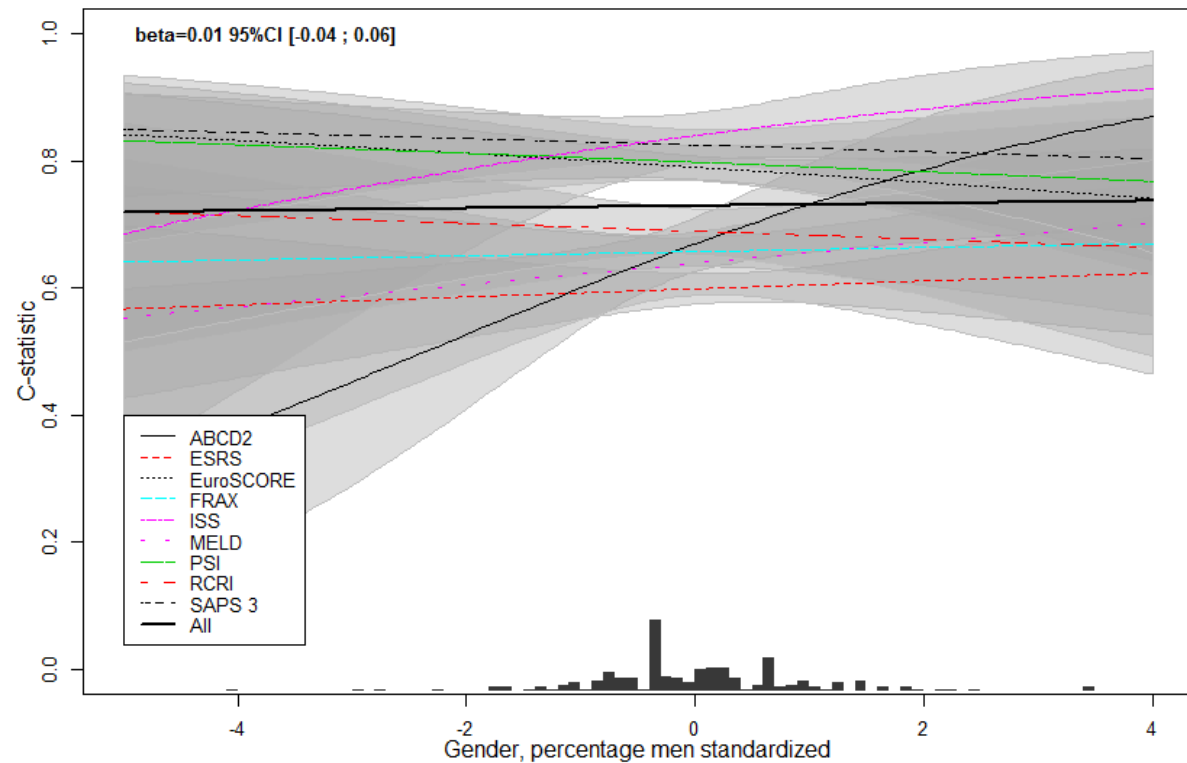
Damen et al. Empirical evidence on the impact of study characteristics and the performance of prediction models: a meta-epidemiological study



Damen et al. Empirical evidence on the impact of study characteristics and the performance of prediction models: a meta-epidemiological study

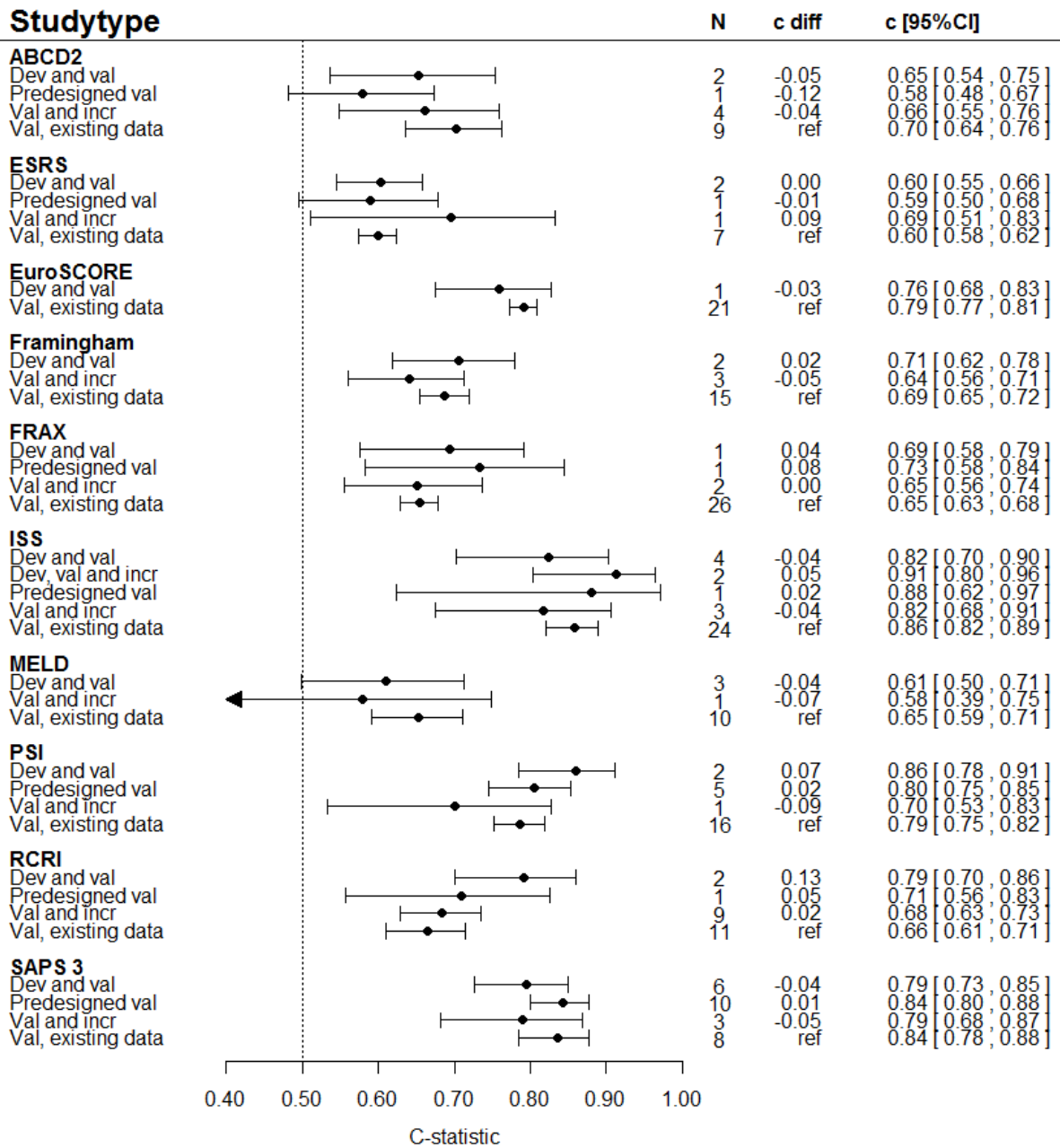


Damen et al. Empirical evidence on the impact of study characteristics and the performance of prediction models: a meta-epidemiological study

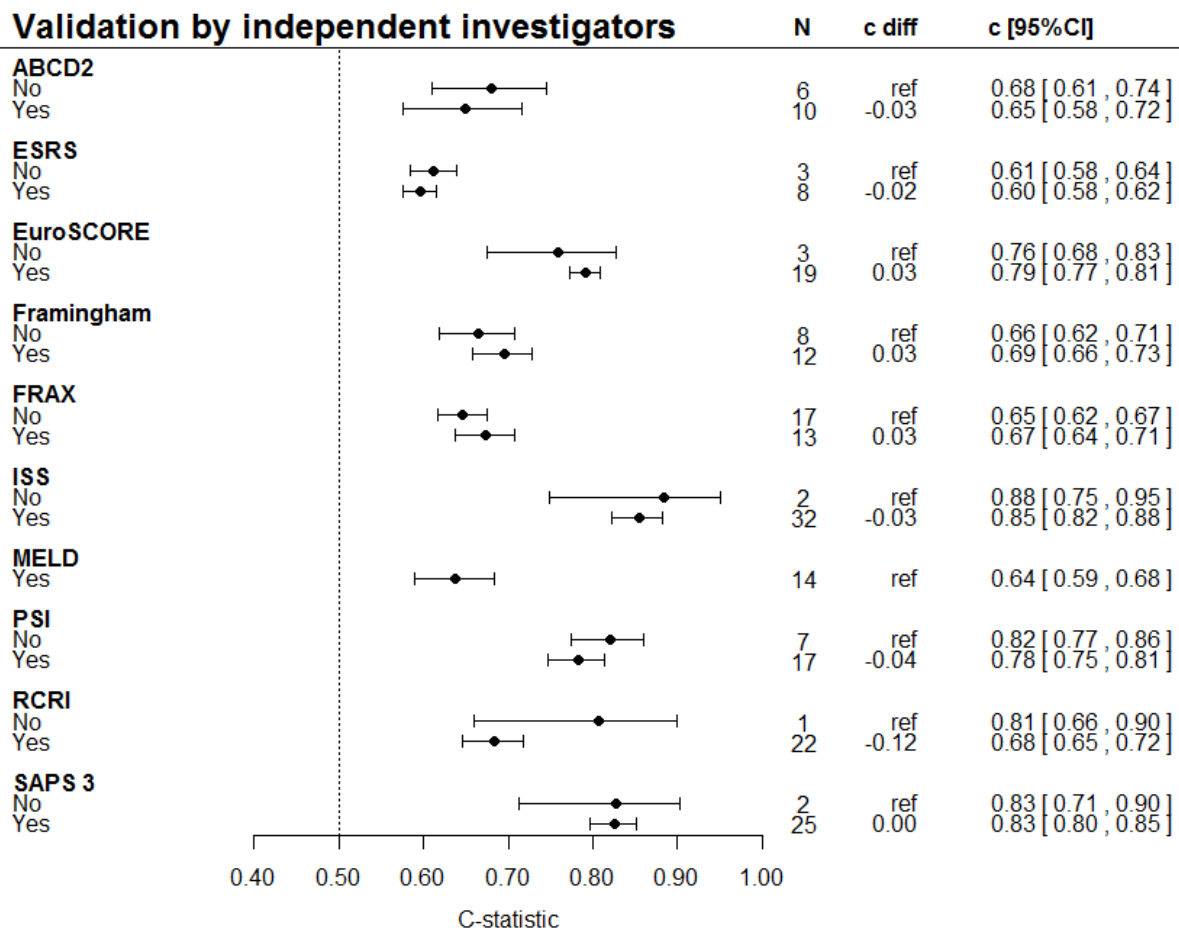


A value of 0 on the x-axis means that the corresponding validation study adopted a common value for that characteristic (i.e. the average value for all validation studies of that specific prediction model), values above 0 mean that values for that characteristic were higher than the average value of all validation studies of that model, whereas values below 0 mean a lower than average value for that characteristic. For example, for the prediction horizon, a value of 0 means that the corresponding validation study adopted the average prediction horizon, values above 0 mean that prediction model performance was assessed for long-term endpoints, whereas values below 0 mean that prediction model performance was assessed for short-term endpoints.

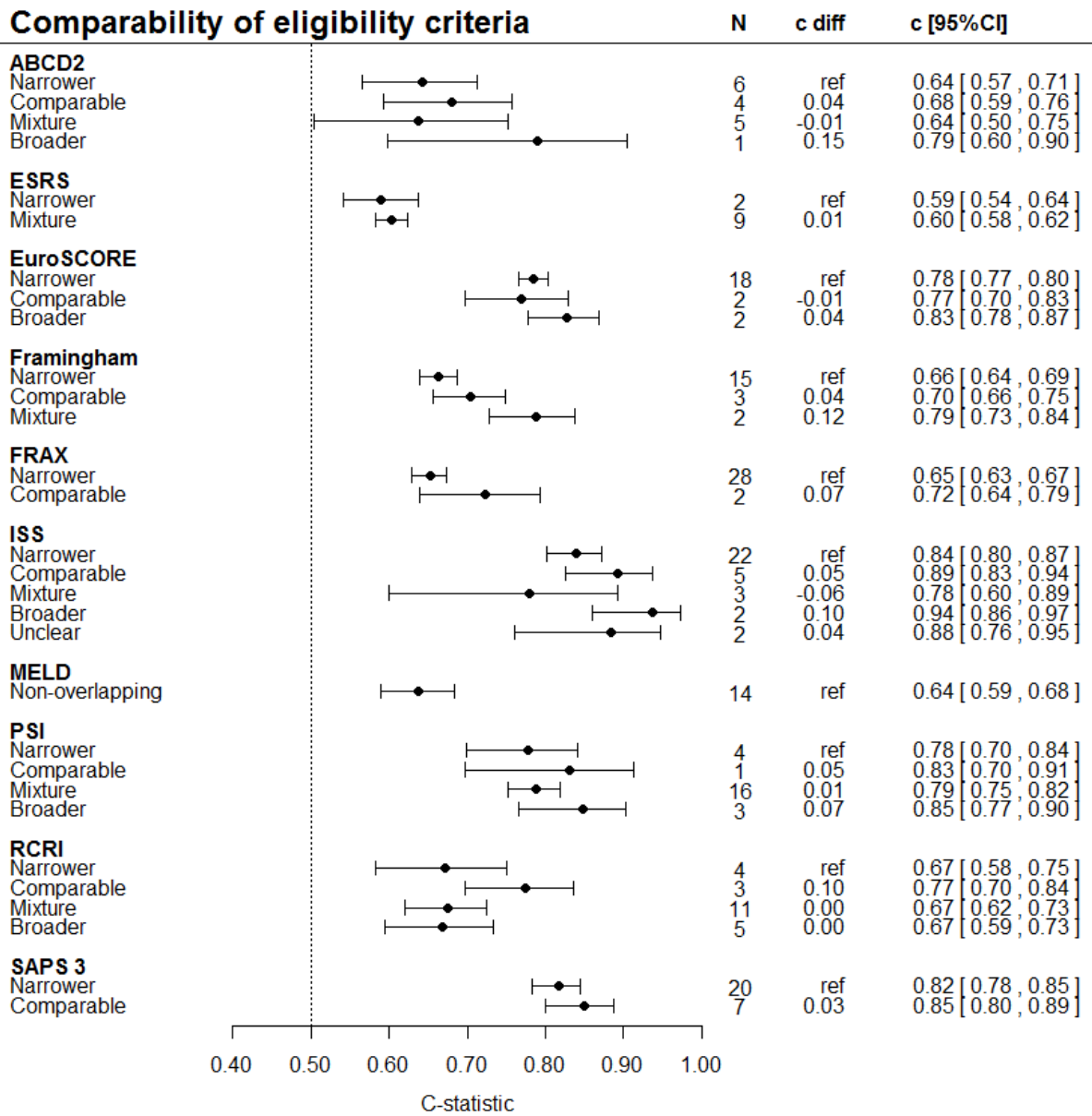
Figure S3: C-statistic in categories of study characteristics within each systematic review



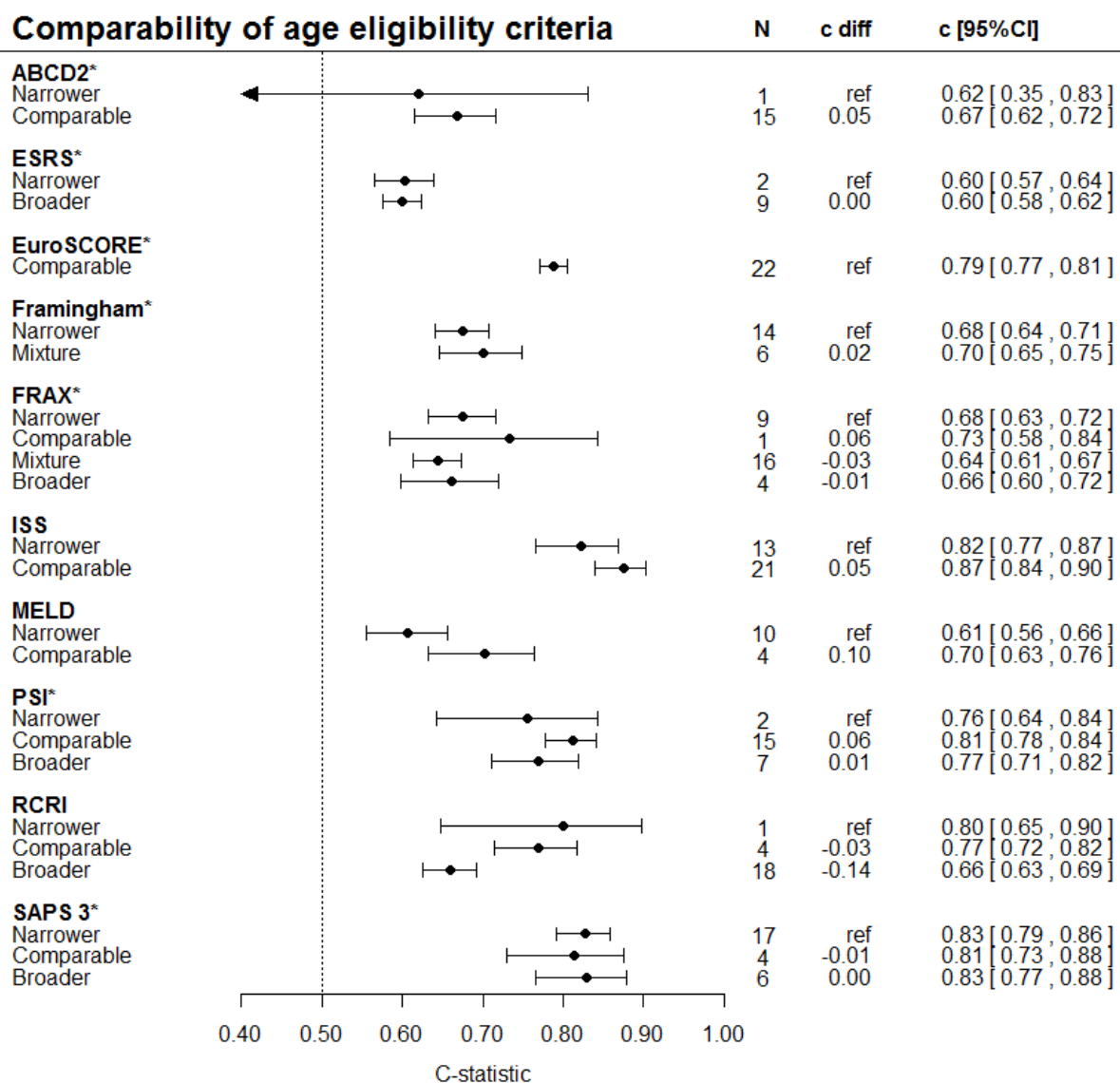
Damen et al. Empirical evidence on the impact of study characteristics and the performance of prediction models: a meta-epidemiological study



Damen et al. Empirical evidence on the impact of study characteristics and the performance of prediction models: a meta-epidemiological study

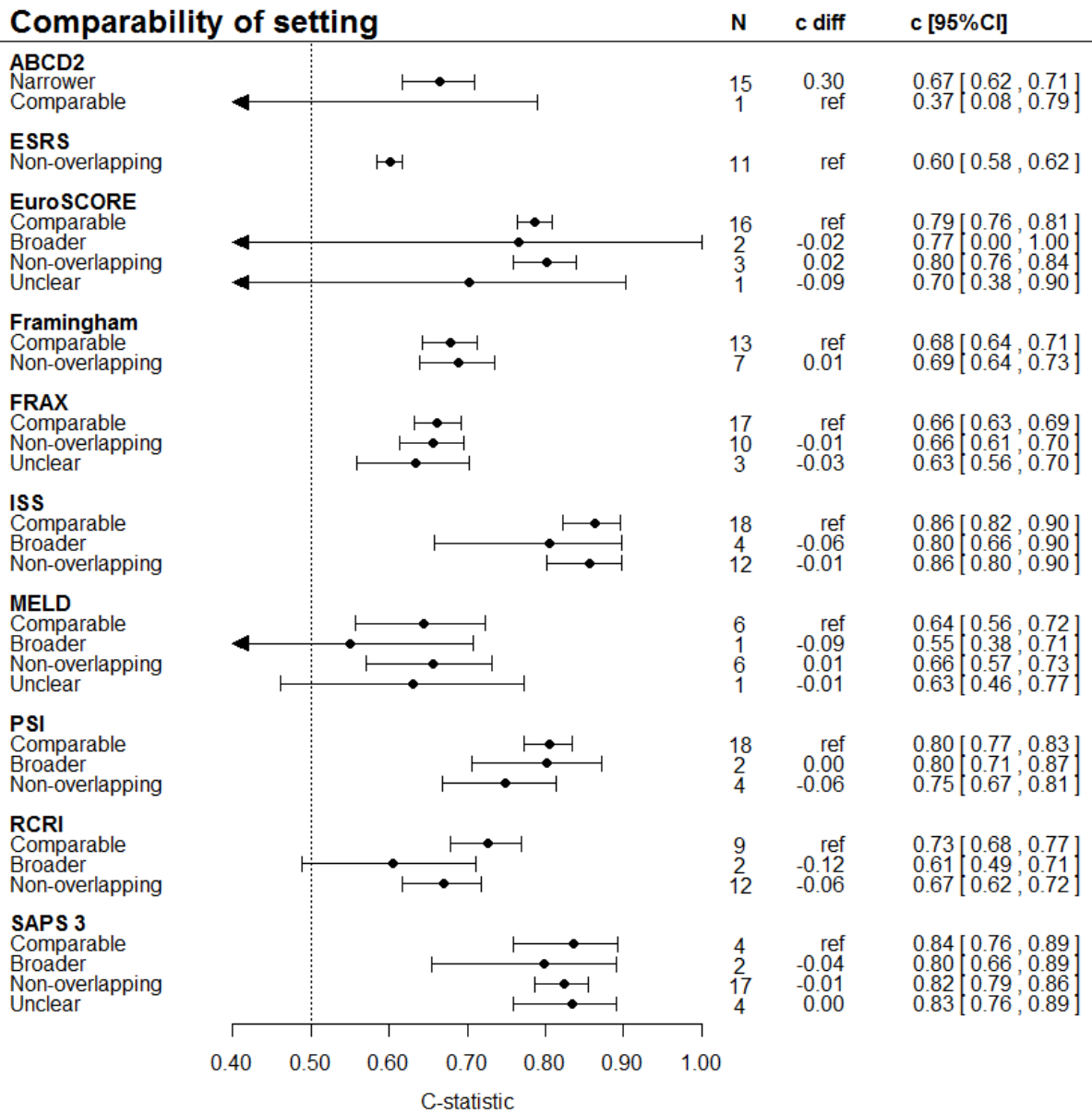


Damen et al. Empirical evidence on the impact of study characteristics and the performance of prediction models: a meta-epidemiological study

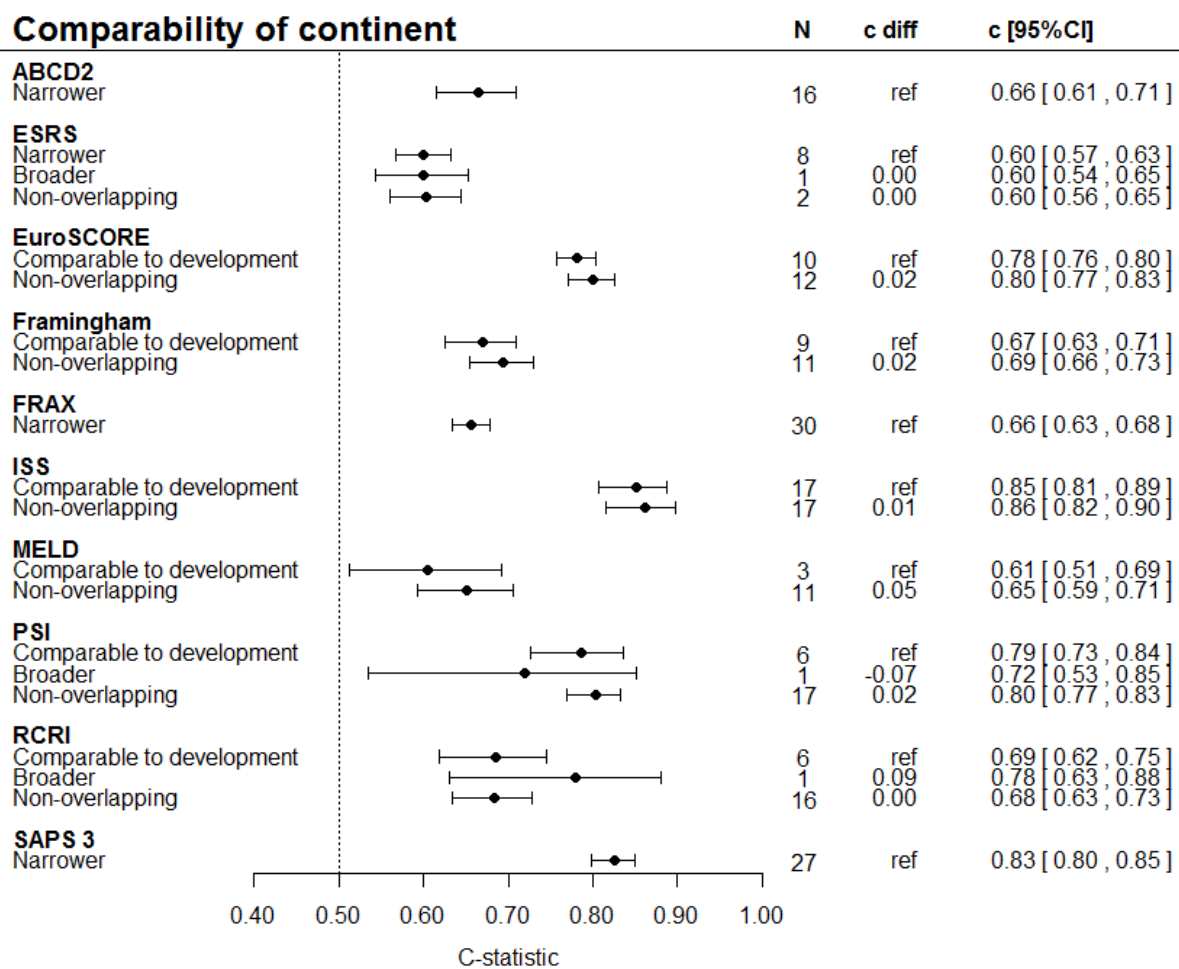


*Models contain age as predictor

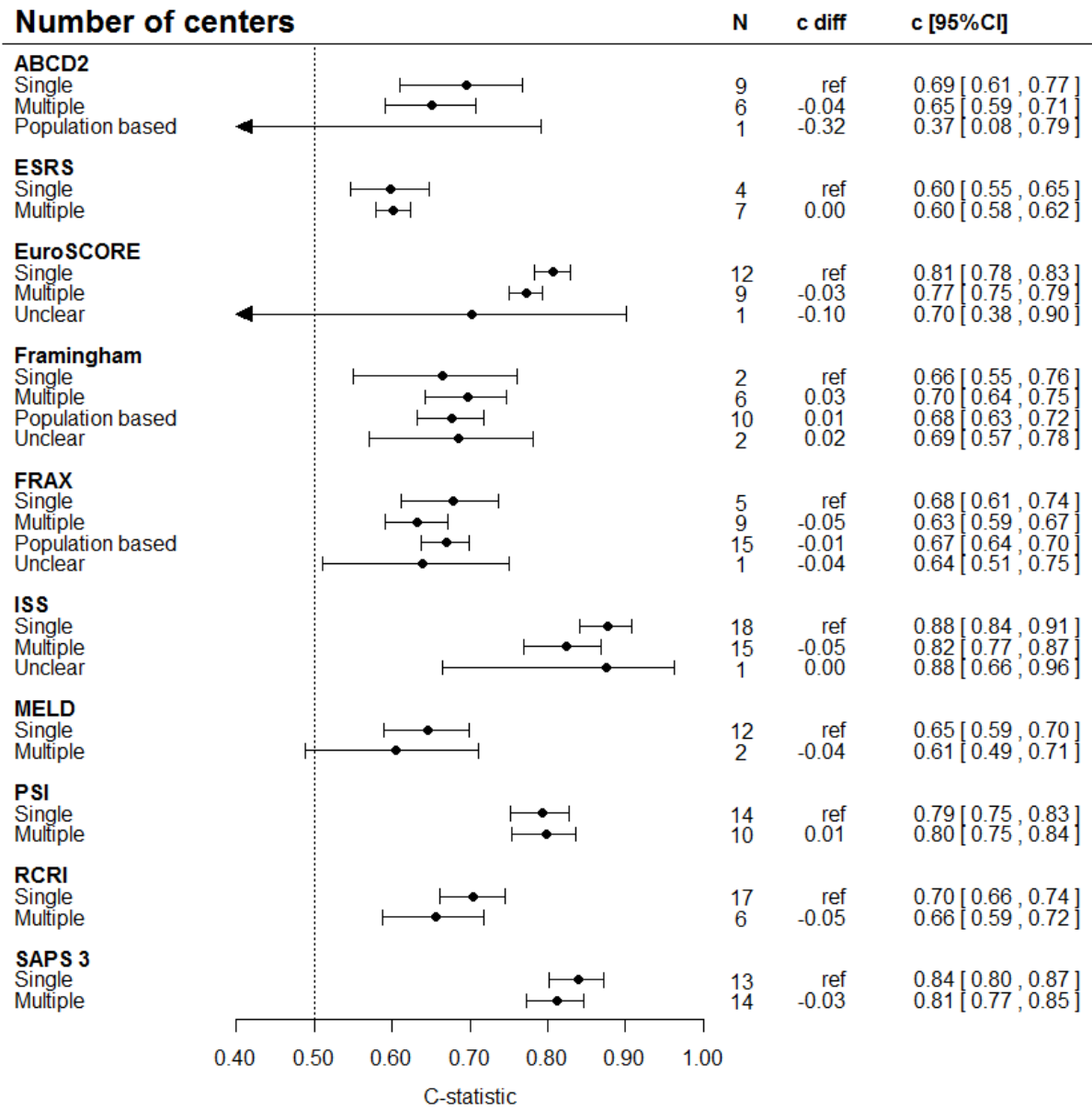
Damen et al. Empirical evidence on the impact of study characteristics and the performance of prediction models: a meta-epidemiological study



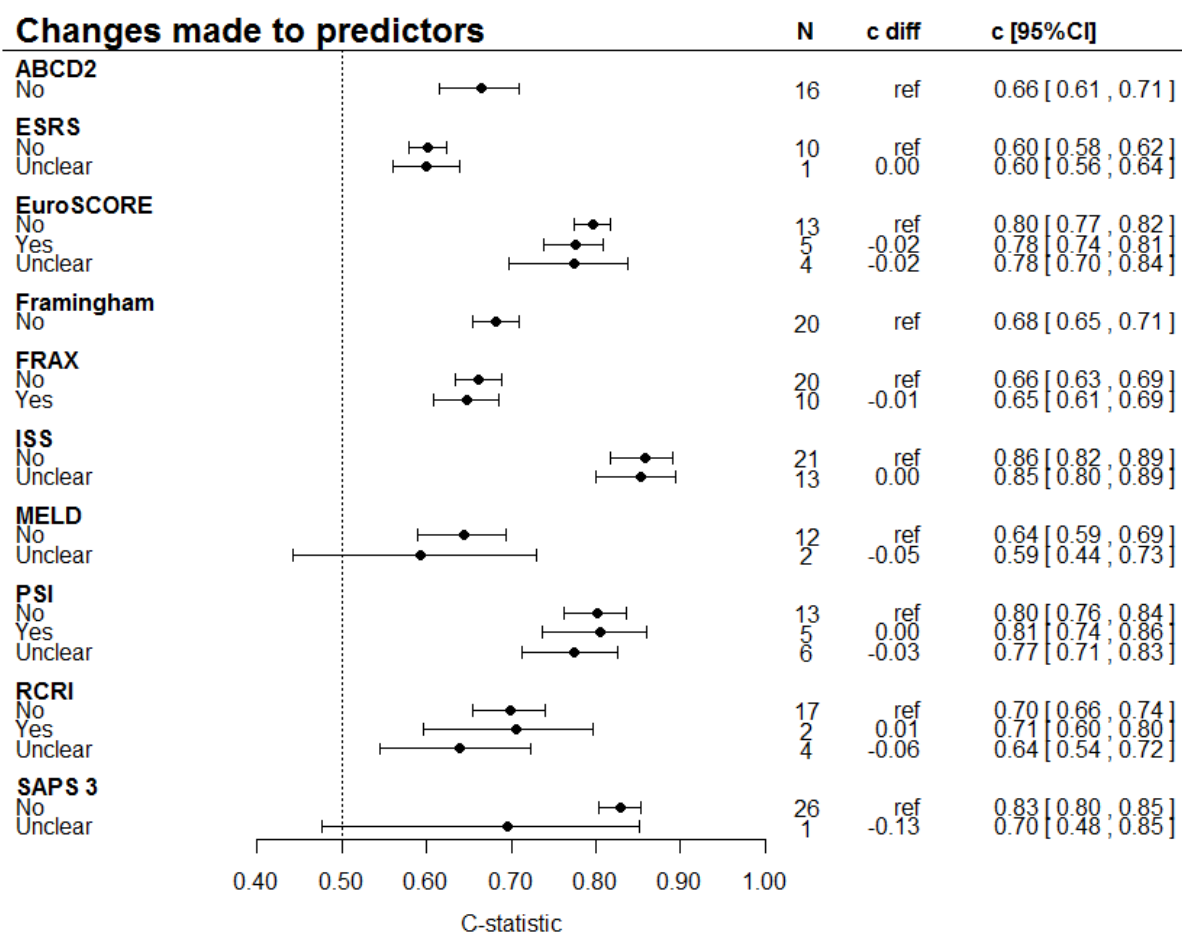
Damen et al. Empirical evidence on the impact of study characteristics and the performance of prediction models: a meta-epidemiological study



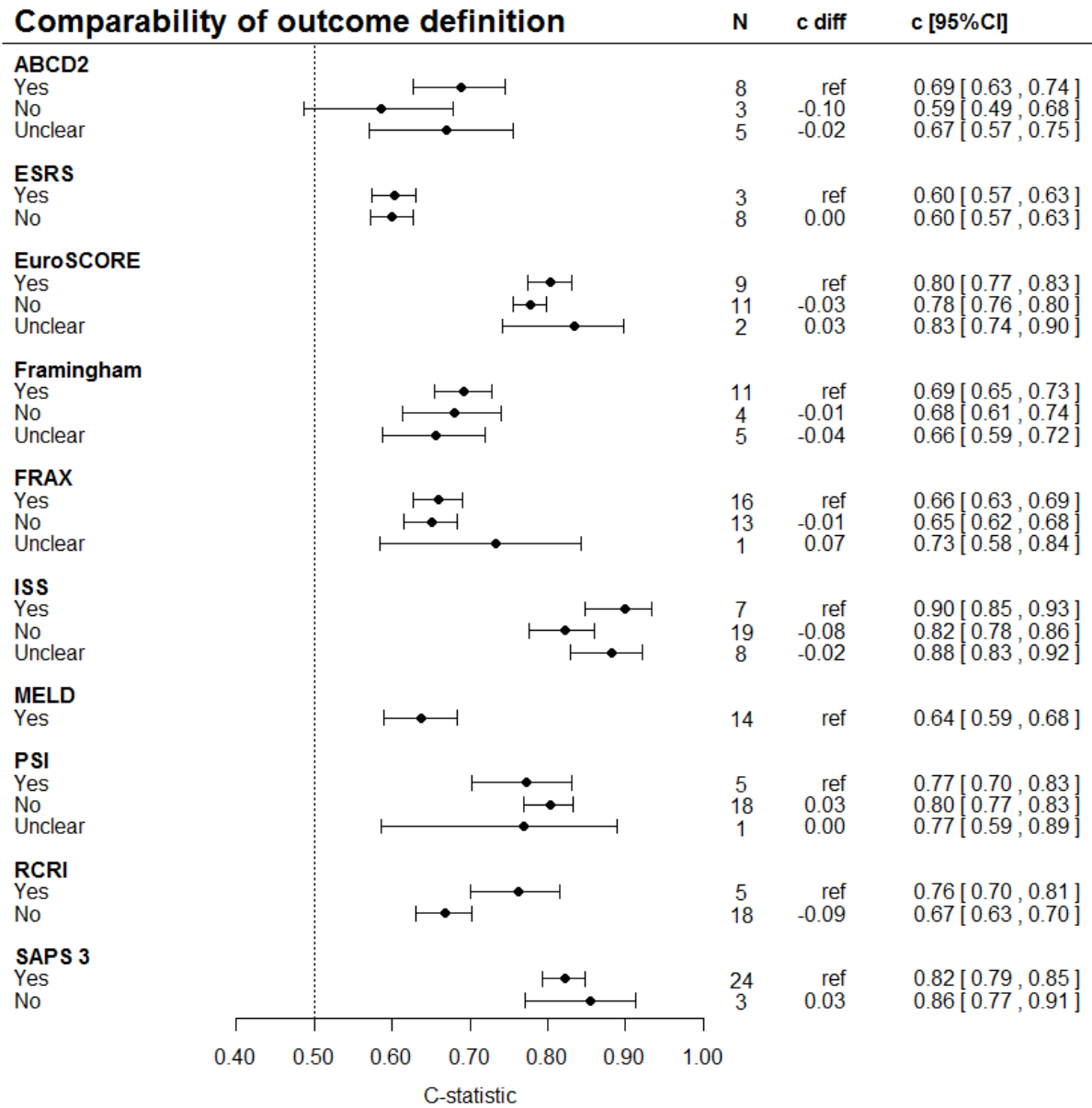
Damen et al. Empirical evidence on the impact of study characteristics and the performance of prediction models: a meta-epidemiological study



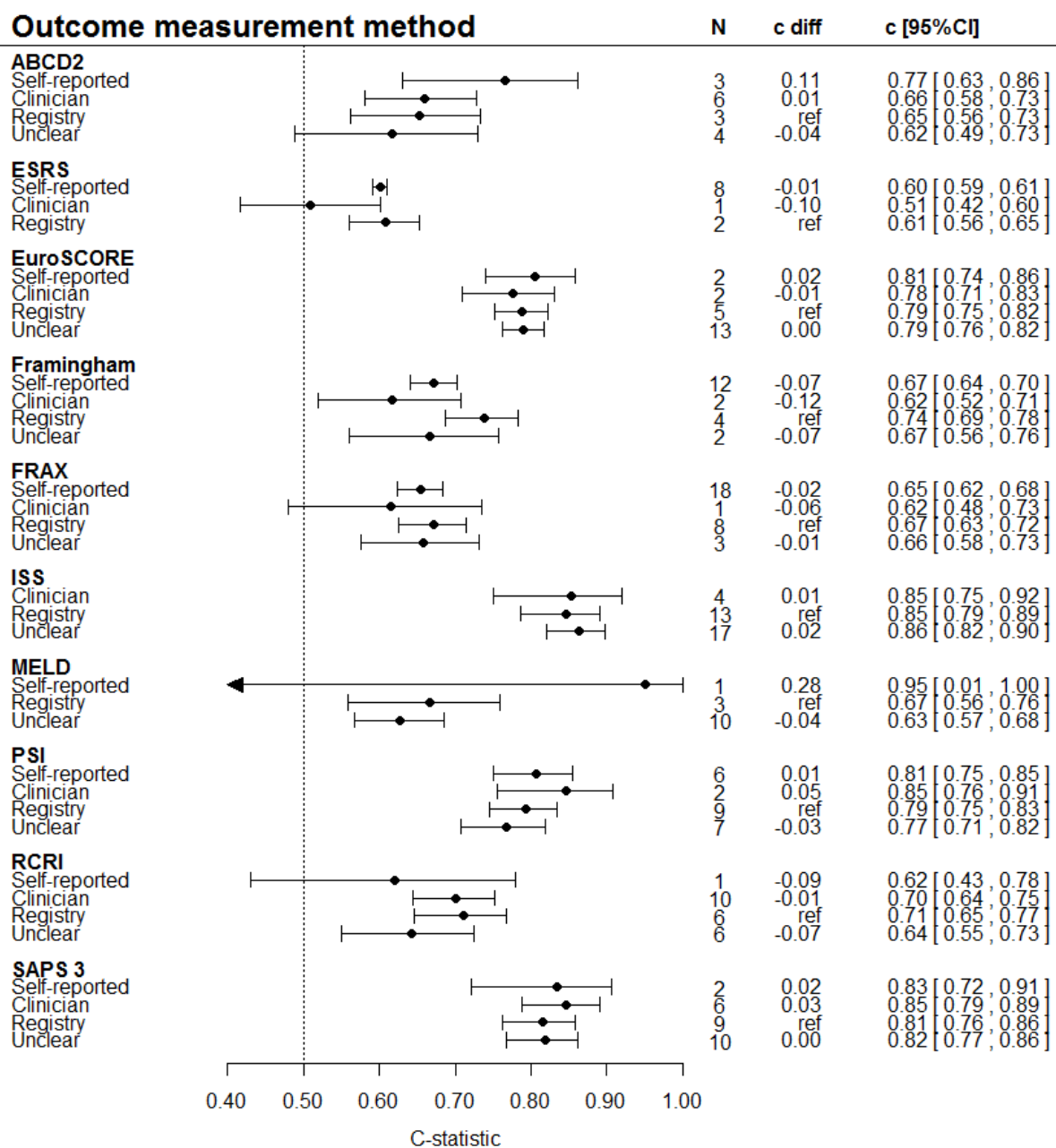
Damen et al. Empirical evidence on the impact of study characteristics and the performance of prediction models: a meta-epidemiological study



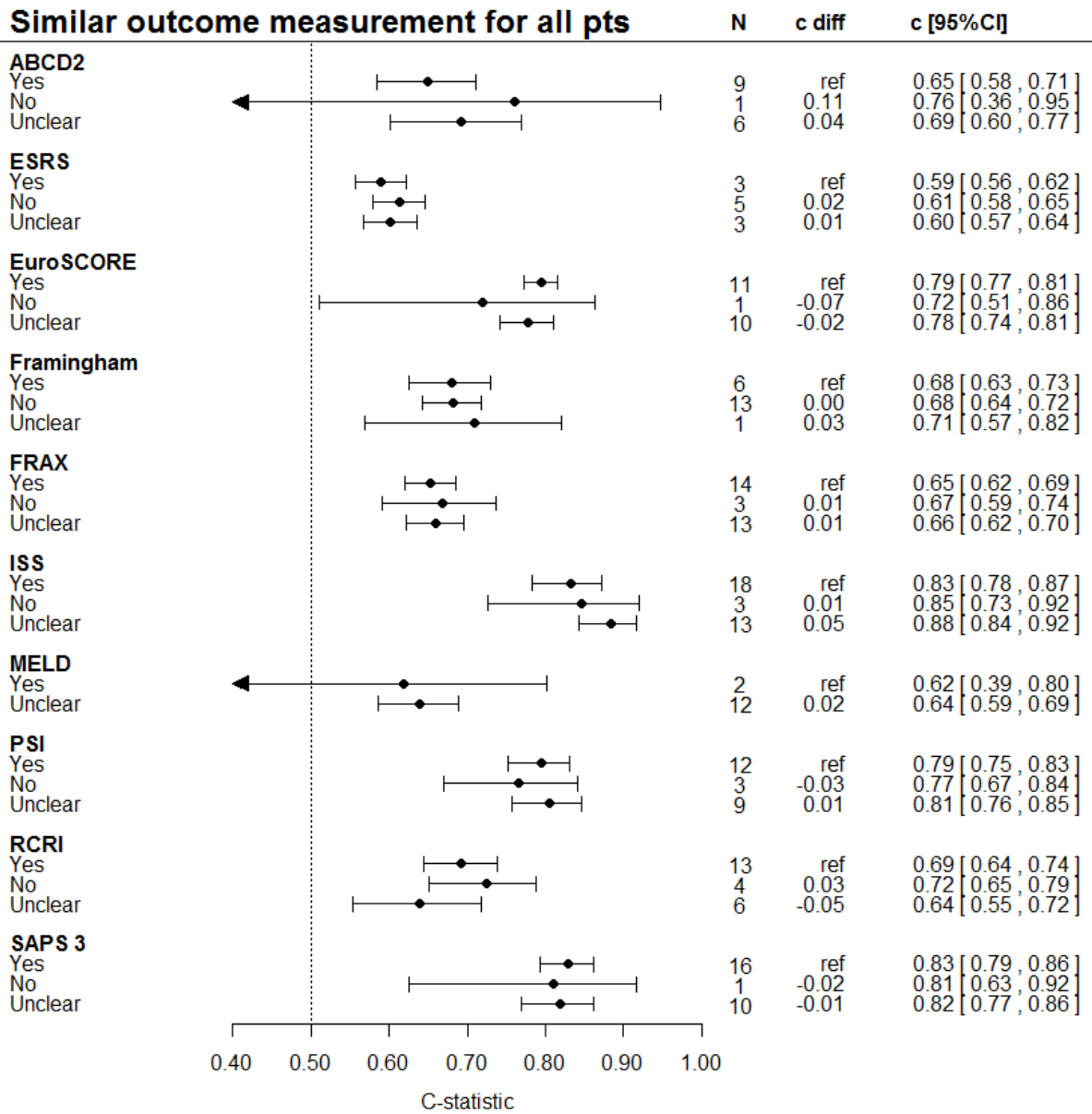
Damen et al. Empirical evidence on the impact of study characteristics and the performance of prediction models: a meta-epidemiological study



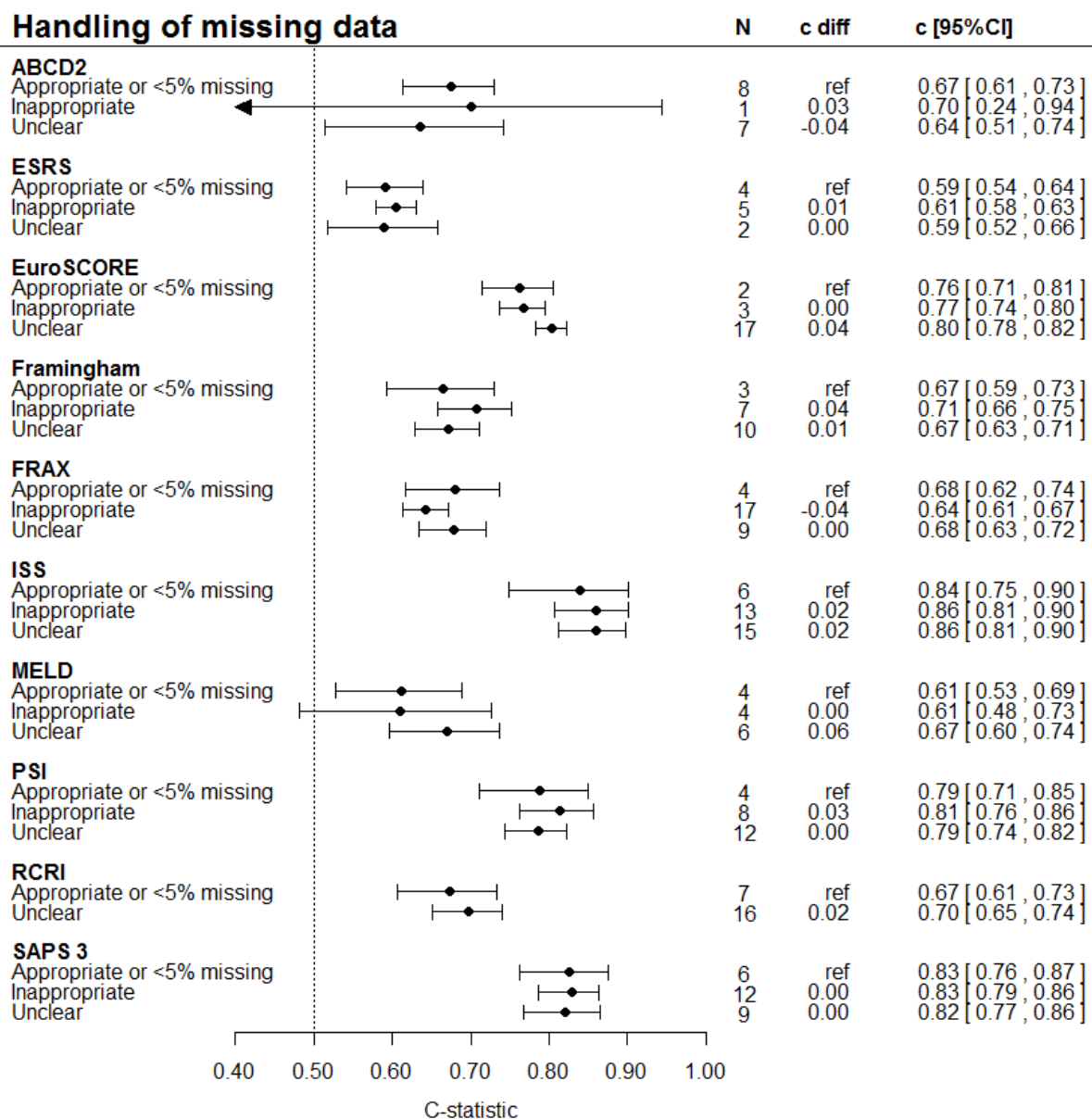
Damen et al. Empirical evidence on the impact of study characteristics and the performance of prediction models: a meta-epidemiological study



Damen et al. Empirical evidence on the impact of study characteristics and the performance of prediction models: a meta-epidemiological study



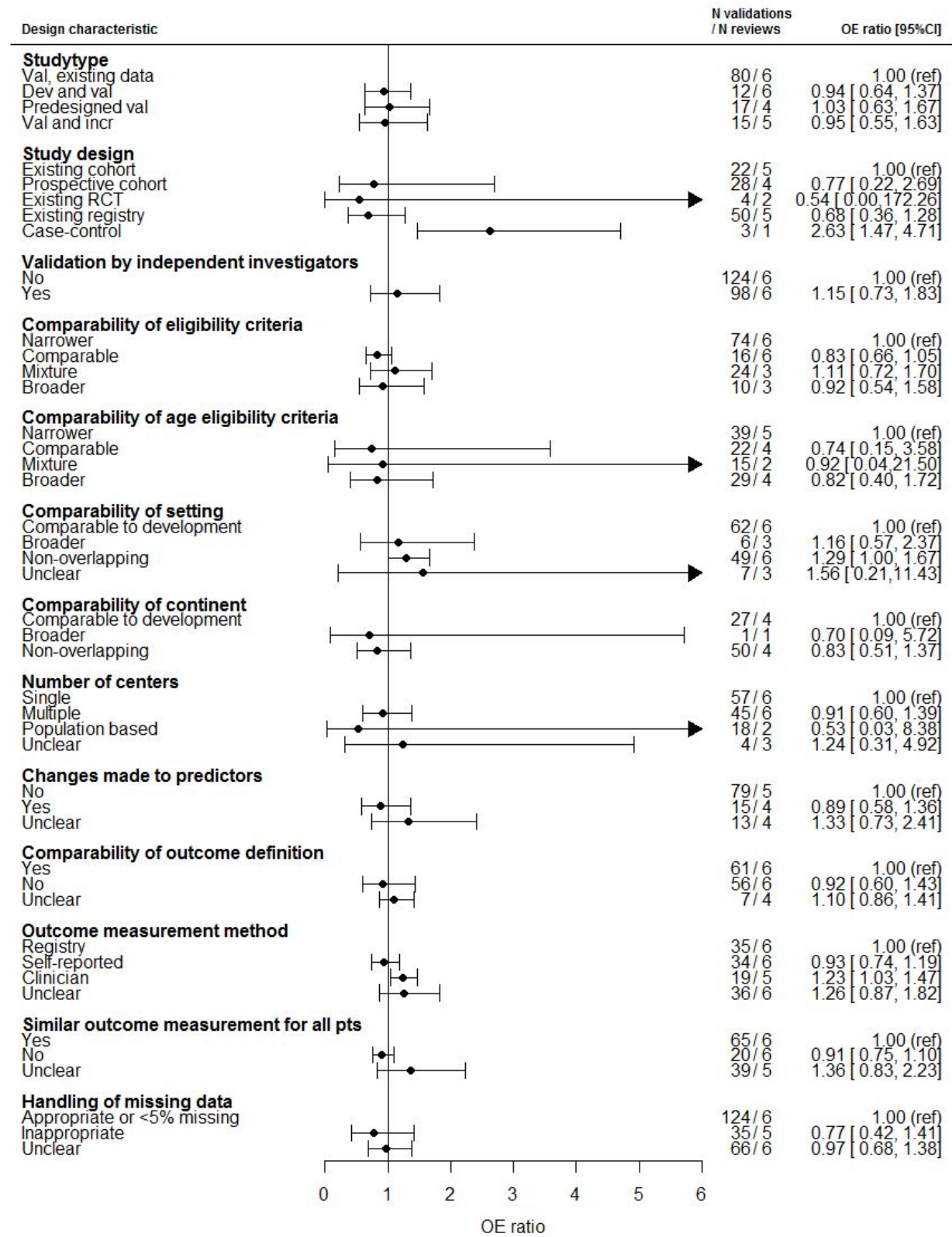
Damen et al. Empirical evidence on the impact of study characteristics and the performance of prediction models: a meta-epidemiological study



C-statistic for categories of study characteristics, pooled using univariable meta-regression analyses per systematic review. N represents the number of external validation studies in a specific category. C diff represents the difference in c-statistic with regard to a reference category (indicated with 'ref'). Dev: development, val: validation, incr: incremental value, pts: patients.

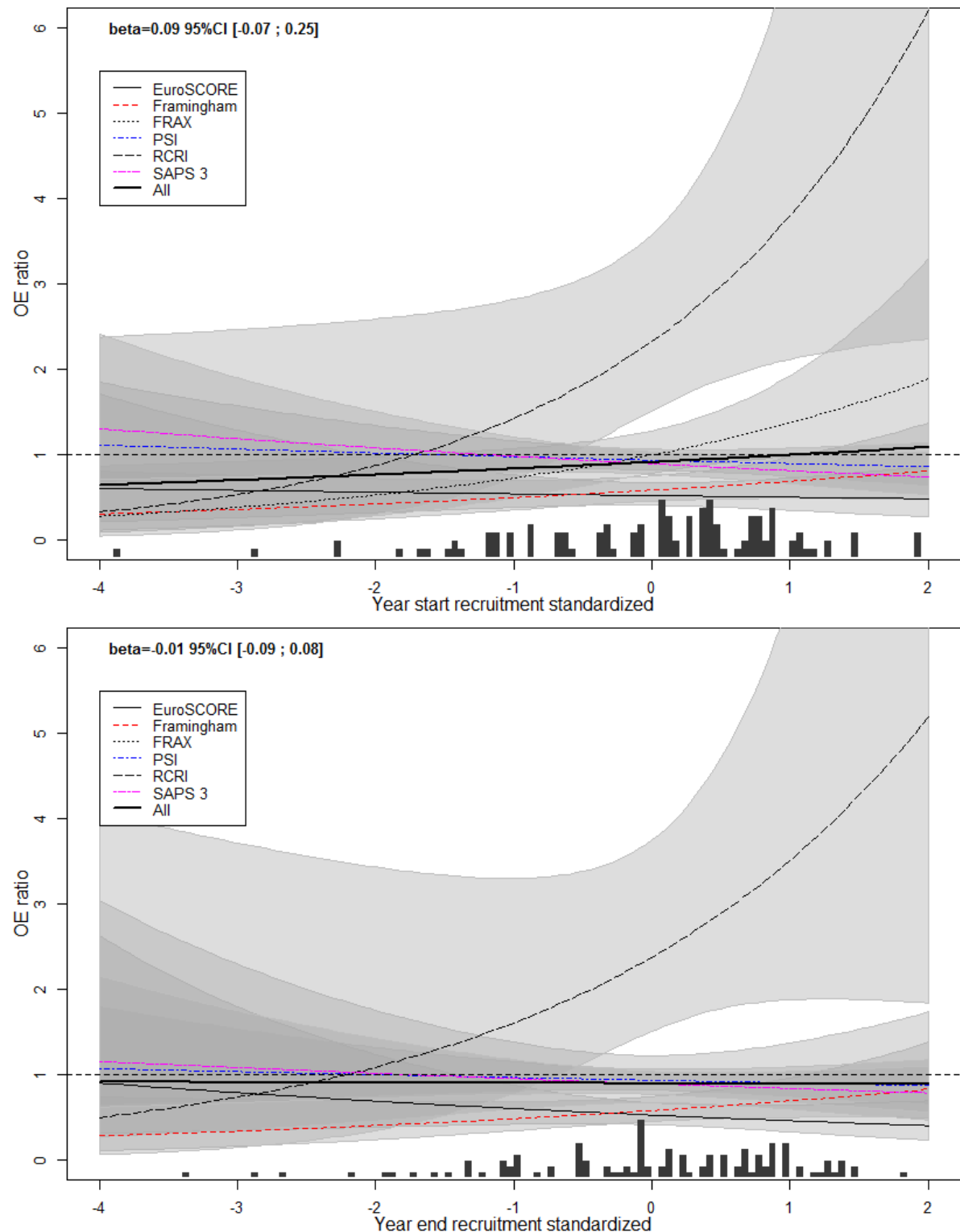
Damen et al. Empirical evidence on the impact of study characteristics and the performance of prediction models: a meta-epidemiological study

Figure S4: Associations between categorical variables and total OE ratio

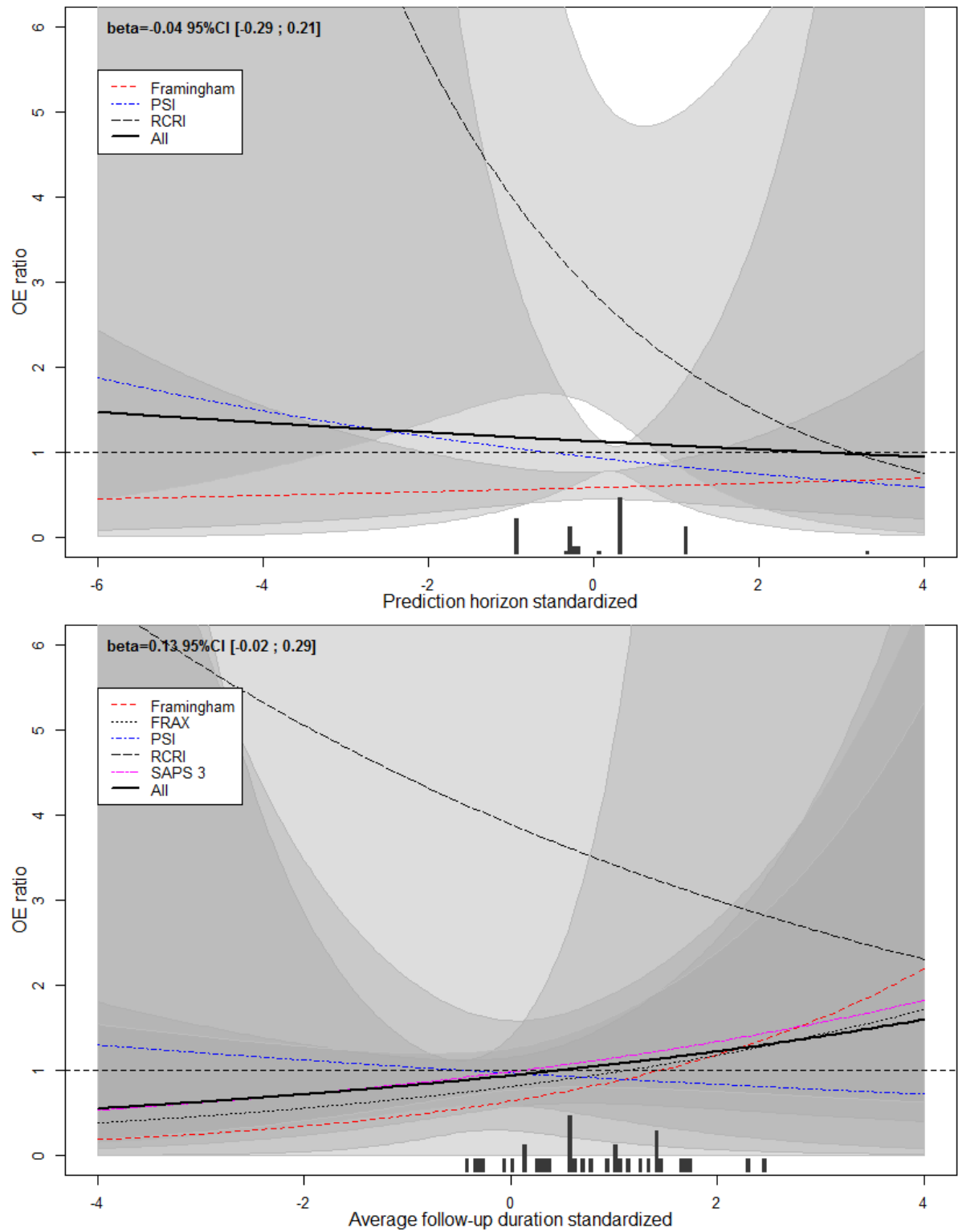


Damen et al. Empirical evidence on the impact of study characteristics and the performance of prediction models: a meta-epidemiological study

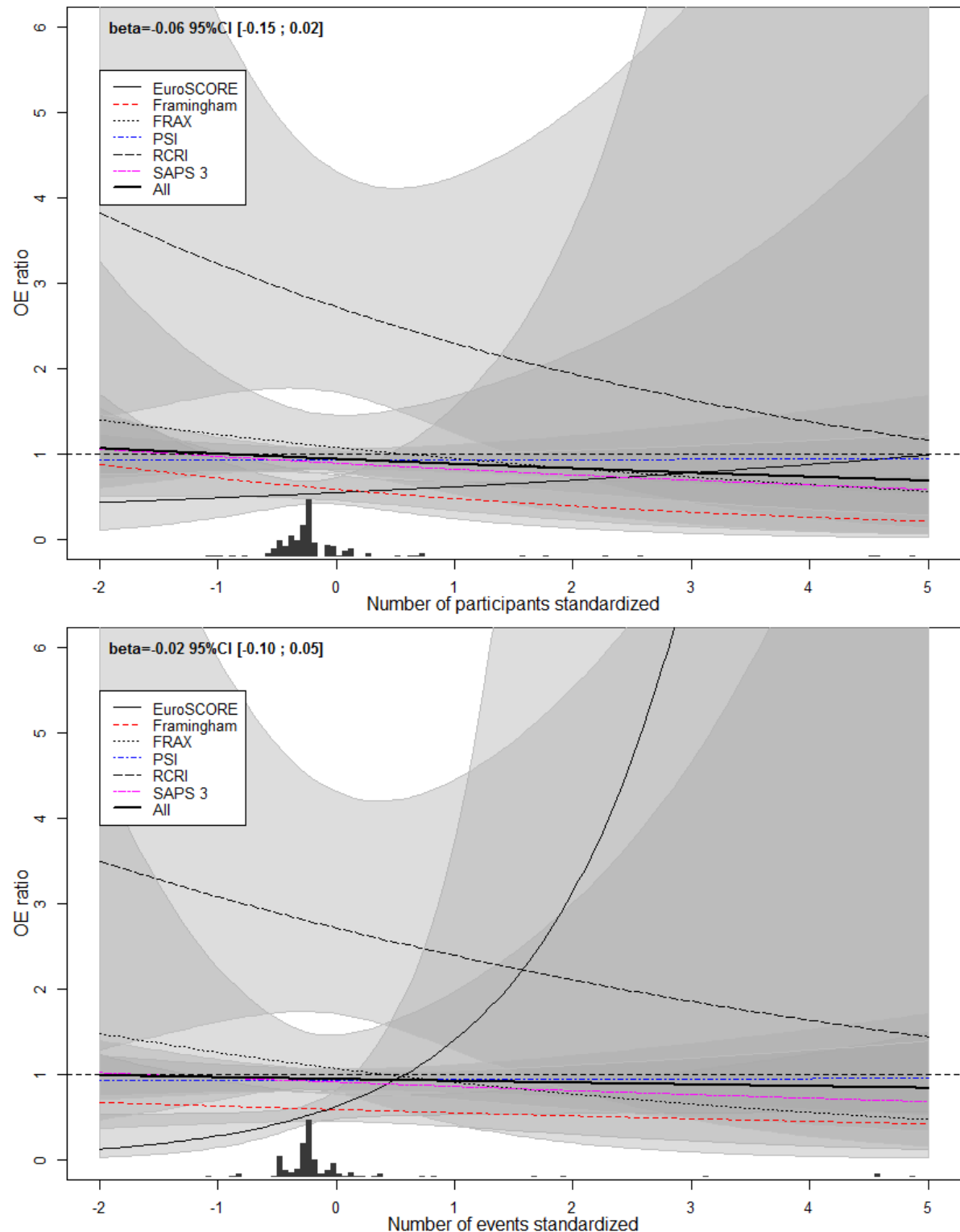
Figure S5: Associations between continuous variables and total OE ratio



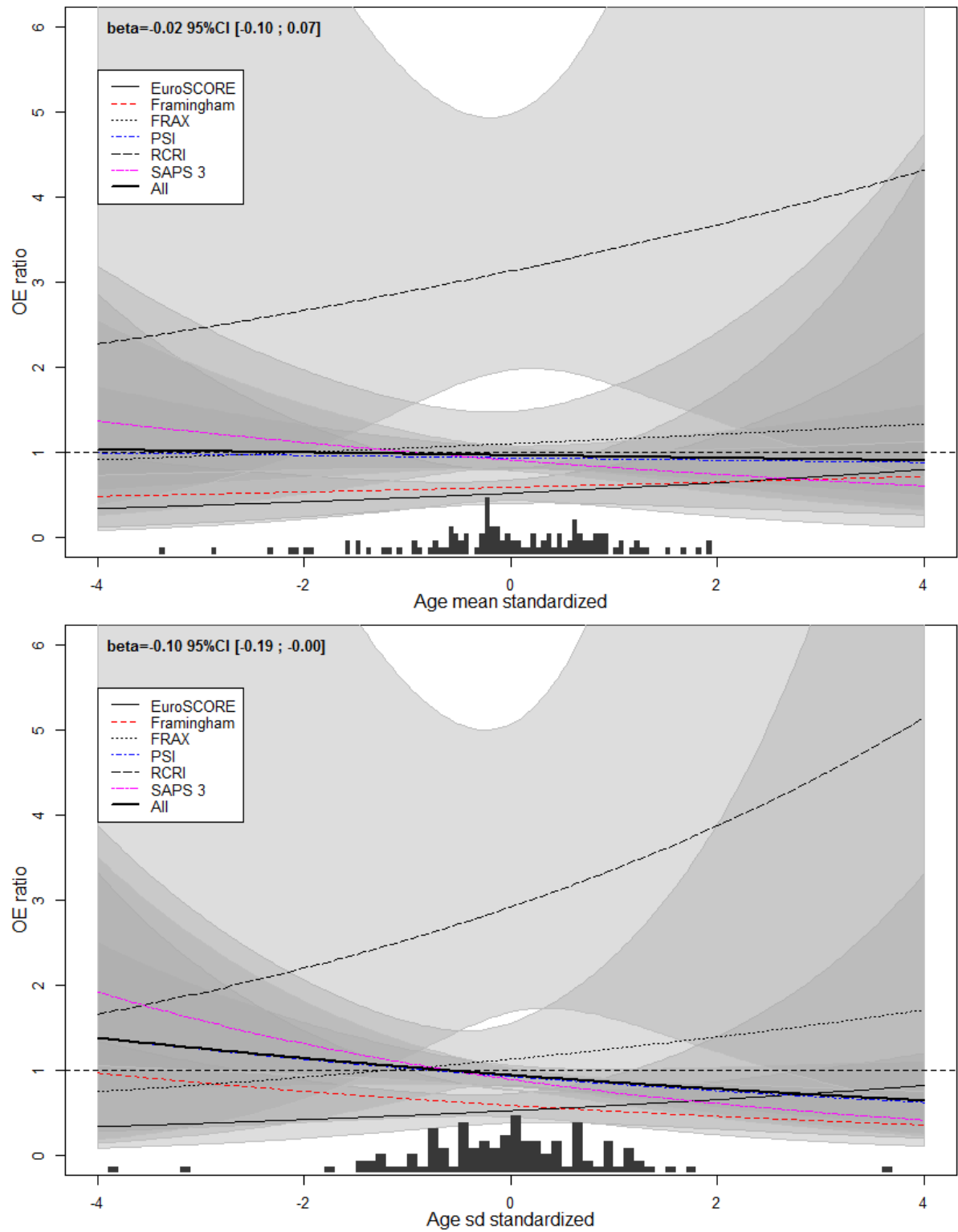
Damen et al. Empirical evidence on the impact of study characteristics and the performance of prediction models: a meta-epidemiological study



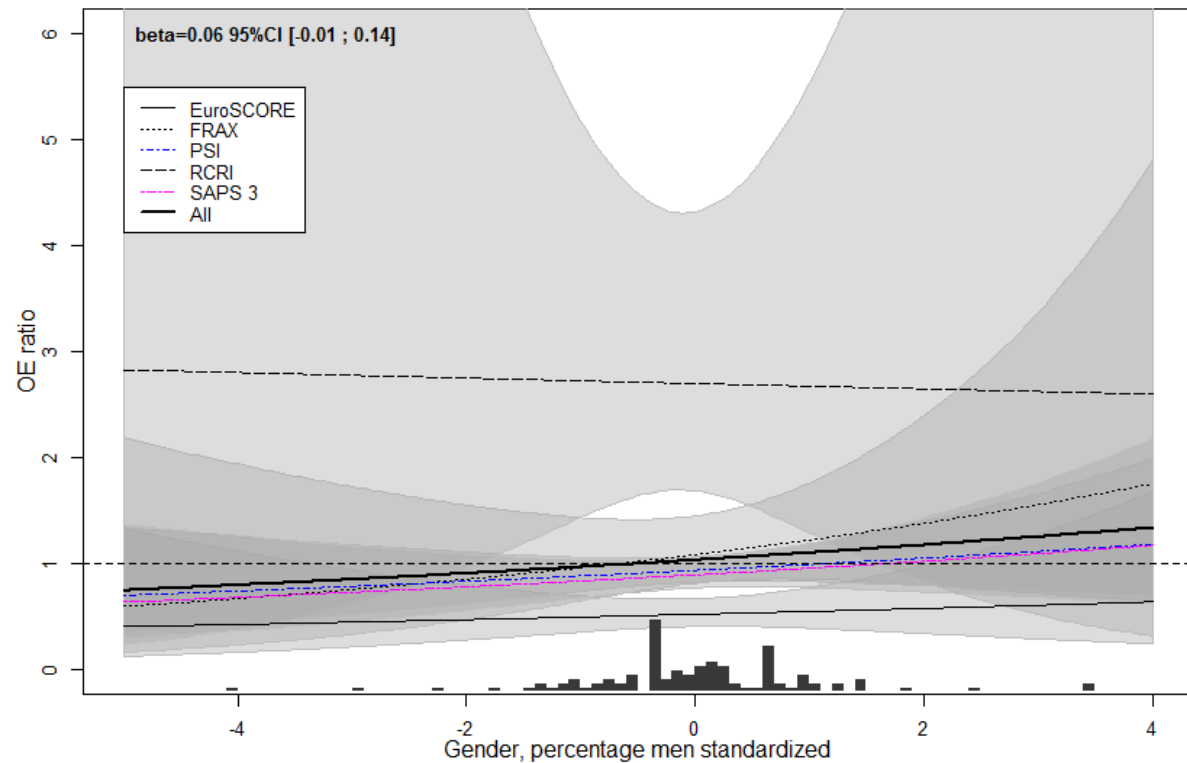
Damen et al. Empirical evidence on the impact of study characteristics and the performance of prediction models: a meta-epidemiological study



Damen et al. Empirical evidence on the impact of study characteristics and the performance of prediction models: a meta-epidemiological study

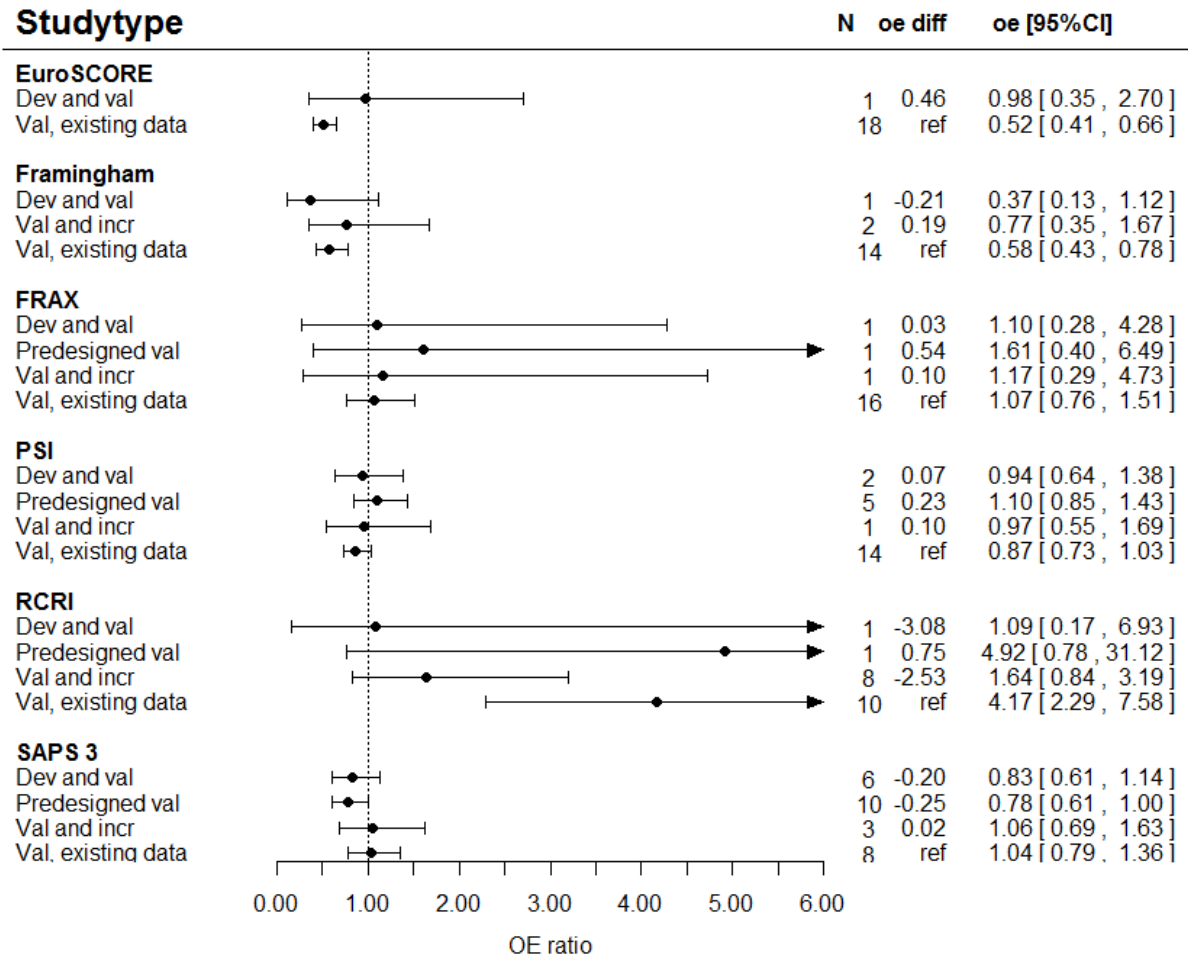


Damen et al. Empirical evidence on the impact of study characteristics and the performance of prediction models: a meta-epidemiological study



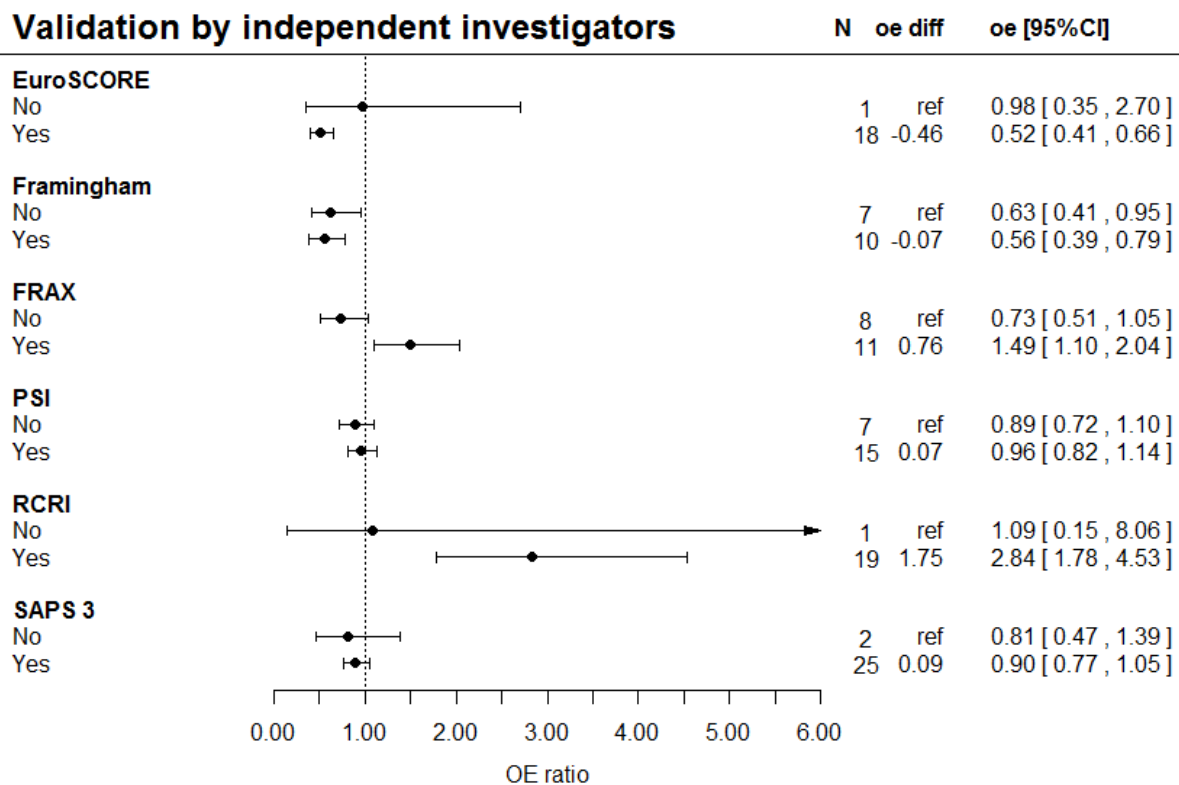
A value of 0 on the x-axis means that the corresponding validation study adopted a common value for that characteristic (i.e. the average value for all validation studies of that specific prediction model), values above 0 mean that values for that characteristic were higher than the average value of all validation studies of that model, whereas values below 0 mean a lower than average value for that characteristic. For example, for the prediction horizon, a value of 0 means that the corresponding validation study adopted the average prediction horizon, values above 0 mean that prediction model performance was assessed for long-term endpoints, whereas values below 0 mean that prediction model performance was assessed for short-term endpoints.

Figure S6: Total OE ratio in categories of study characteristics within each systematic review

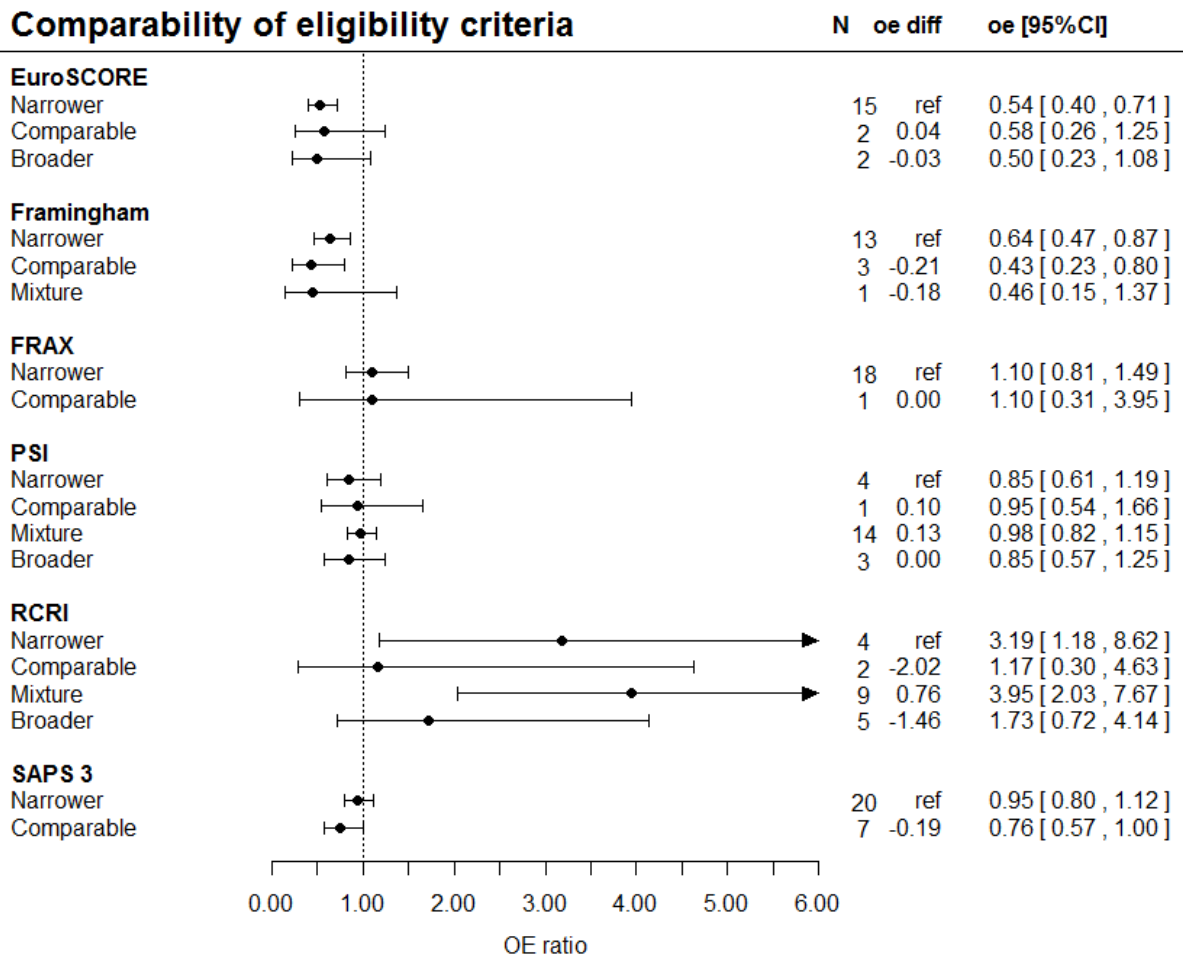


Damen et al. Empirical evidence on the impact of study characteristics and the performance of prediction models: a meta-epidemiological study

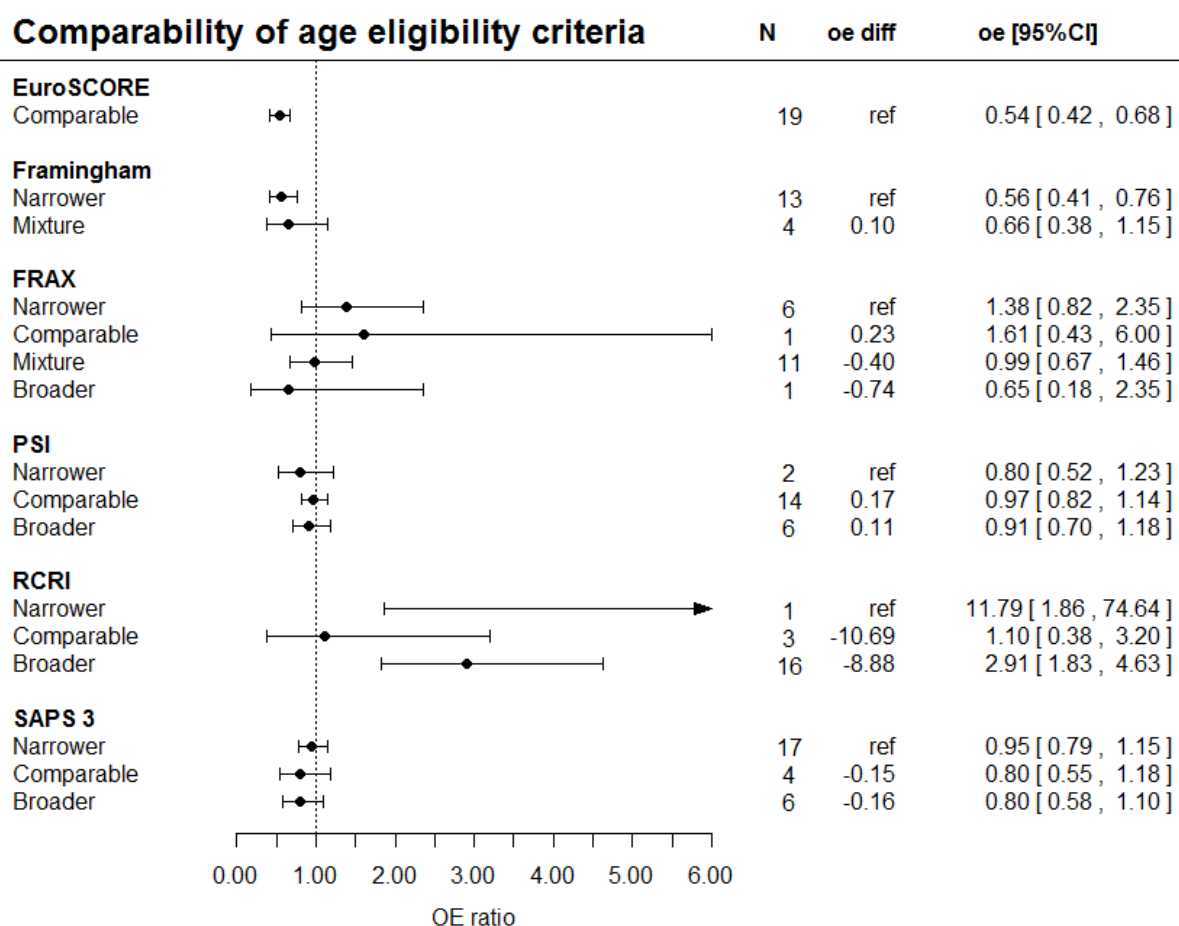
Validation by independent investigators



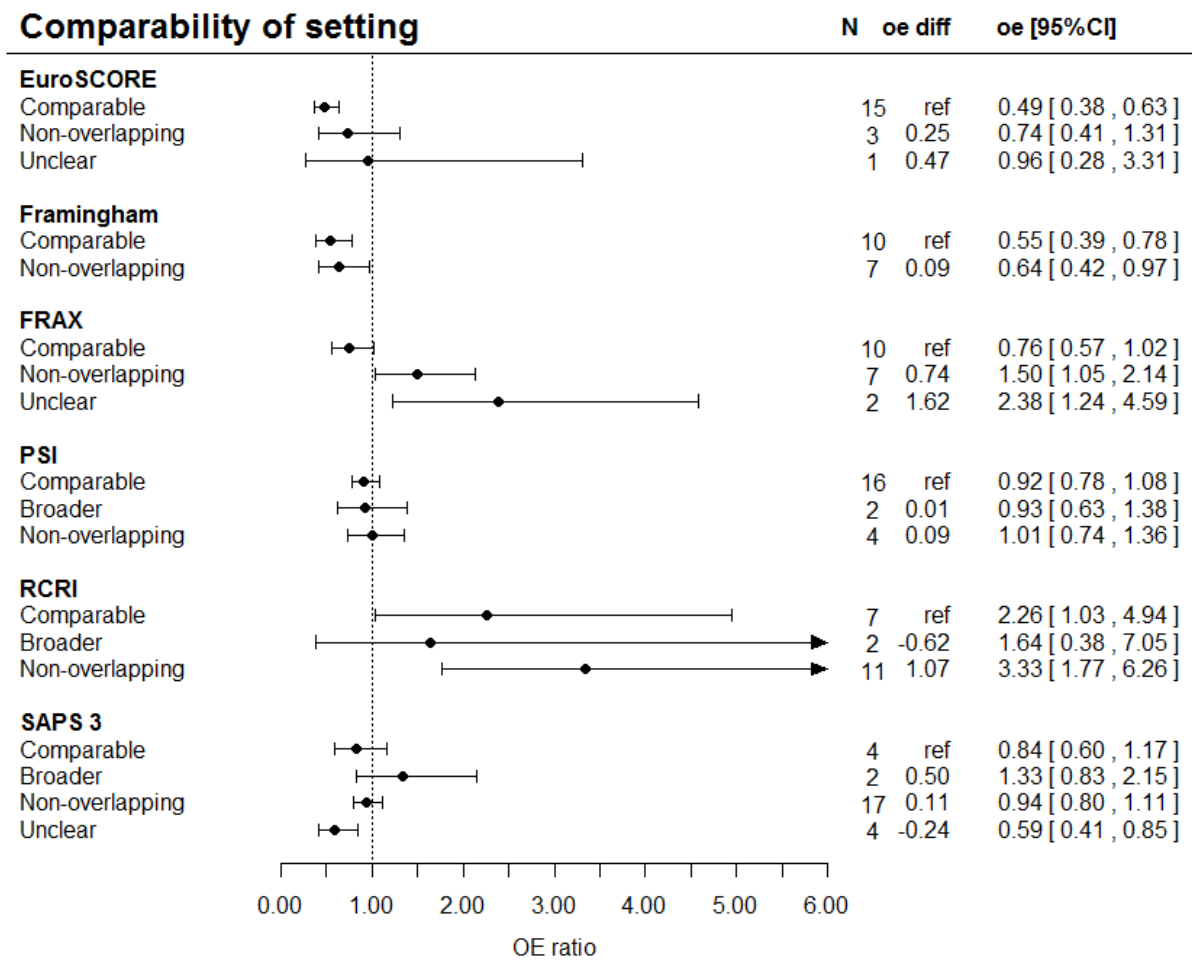
Damen et al. Empirical evidence on the impact of study characteristics and the performance of prediction models: a meta-epidemiological study



Damen et al. Empirical evidence on the impact of study characteristics and the performance of prediction models: a meta-epidemiological study

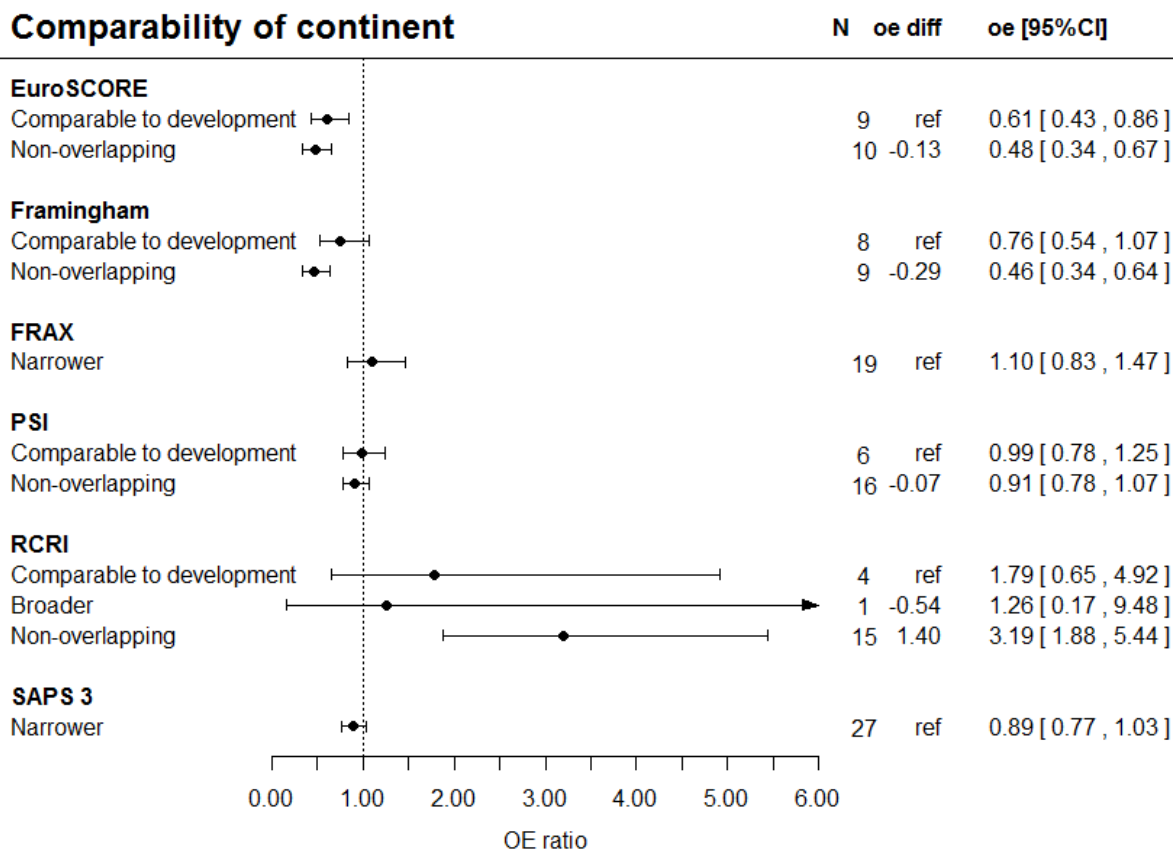


Damen et al. Empirical evidence on the impact of study characteristics and the performance of prediction models: a meta-epidemiological study

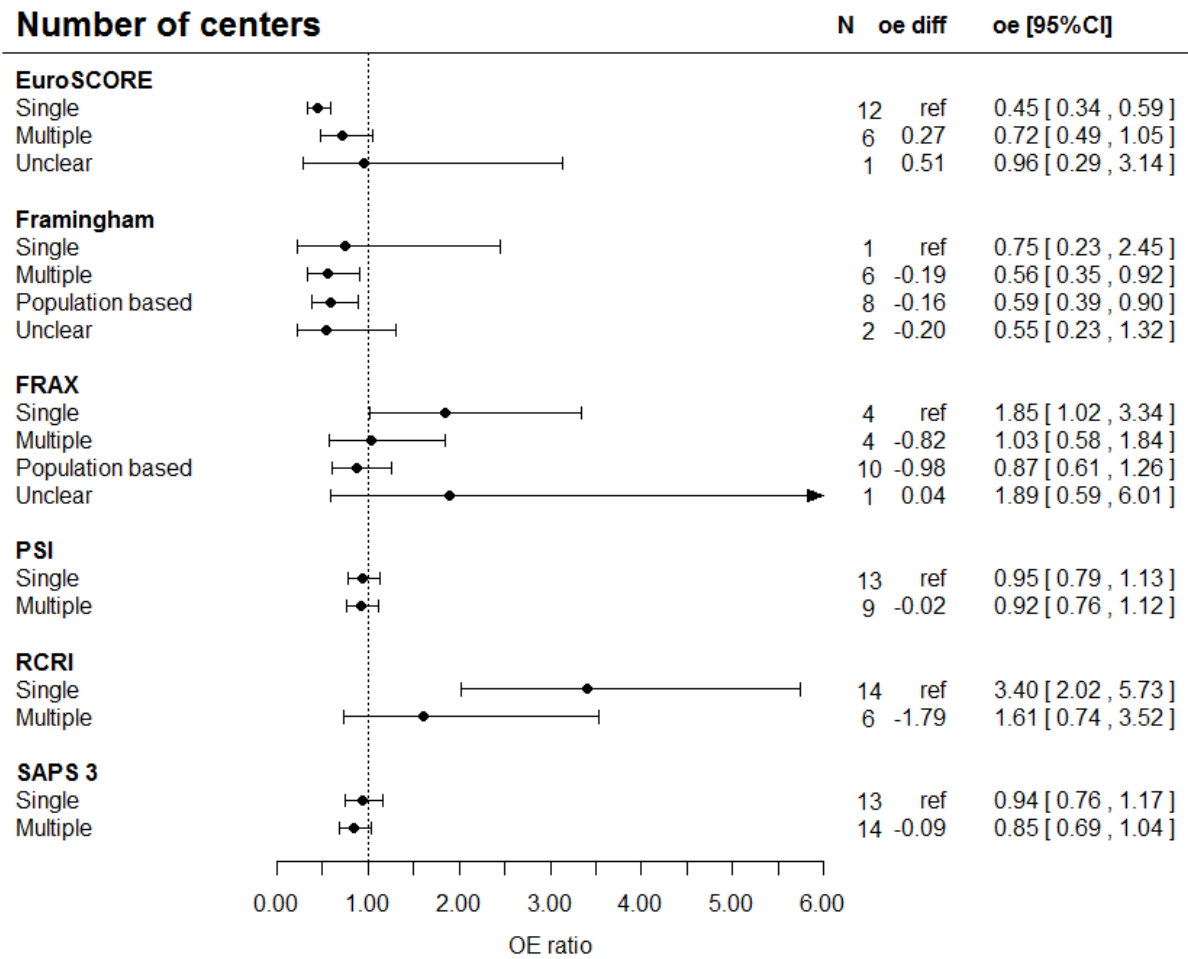


Damen et al. Empirical evidence on the impact of study characteristics and the performance of prediction models: a meta-epidemiological study

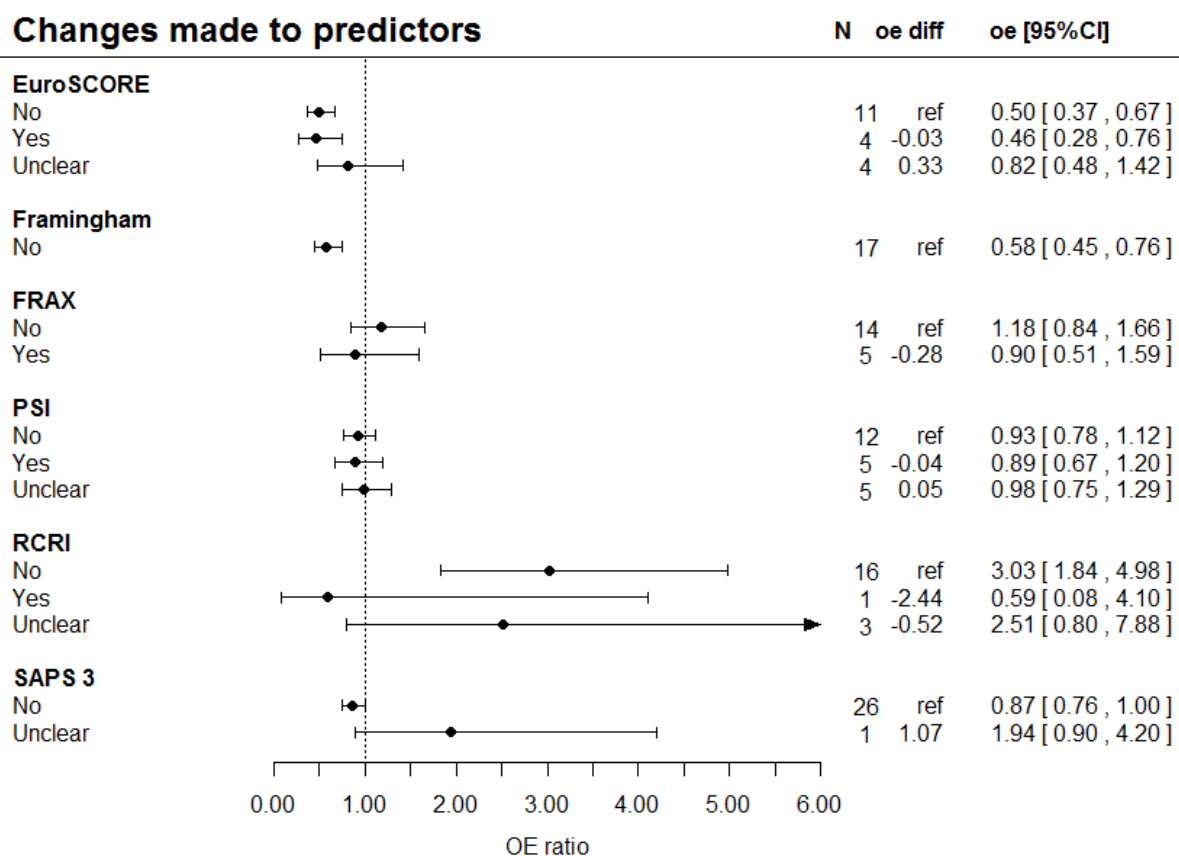
Comparability of continent



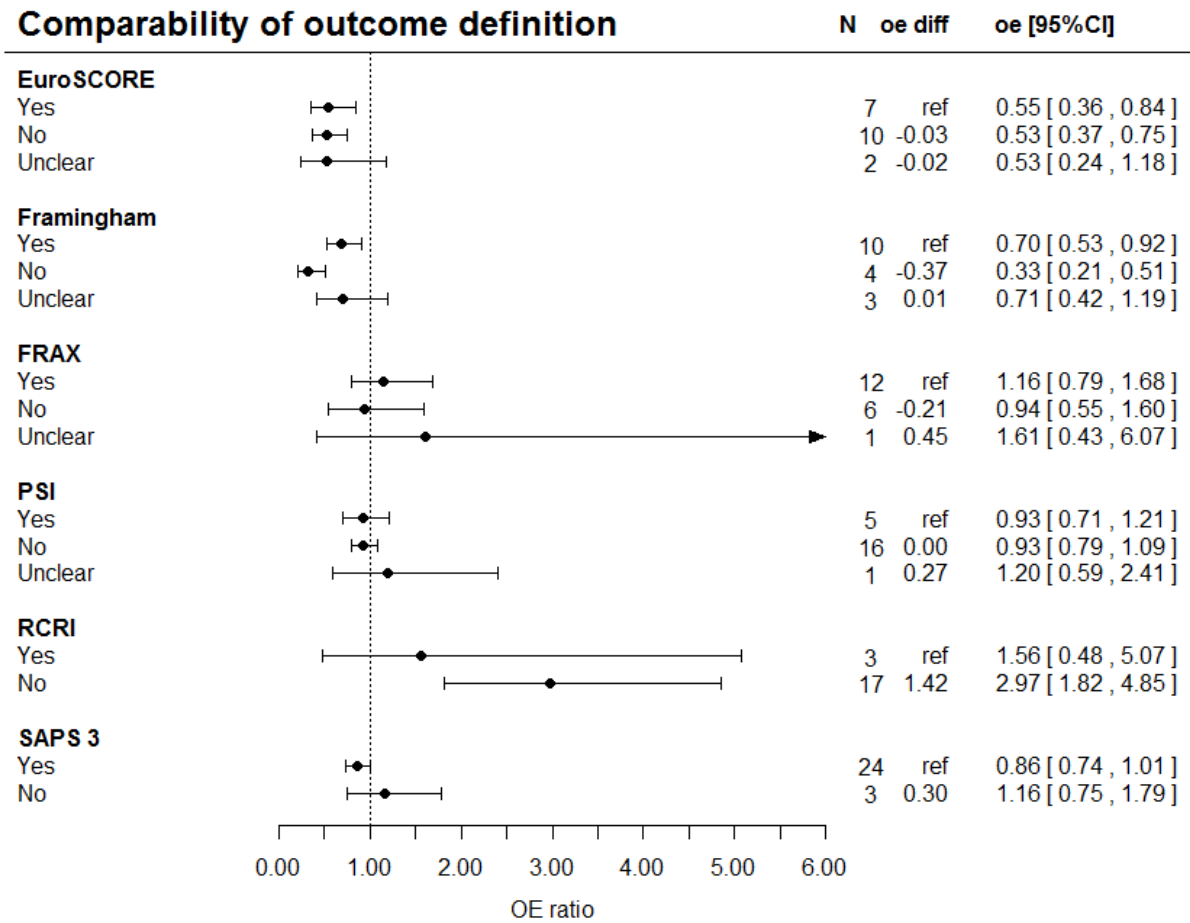
Damen et al. Empirical evidence on the impact of study characteristics and the performance of prediction models: a meta-epidemiological study



Damen et al. Empirical evidence on the impact of study characteristics and the performance of prediction models: a meta-epidemiological study

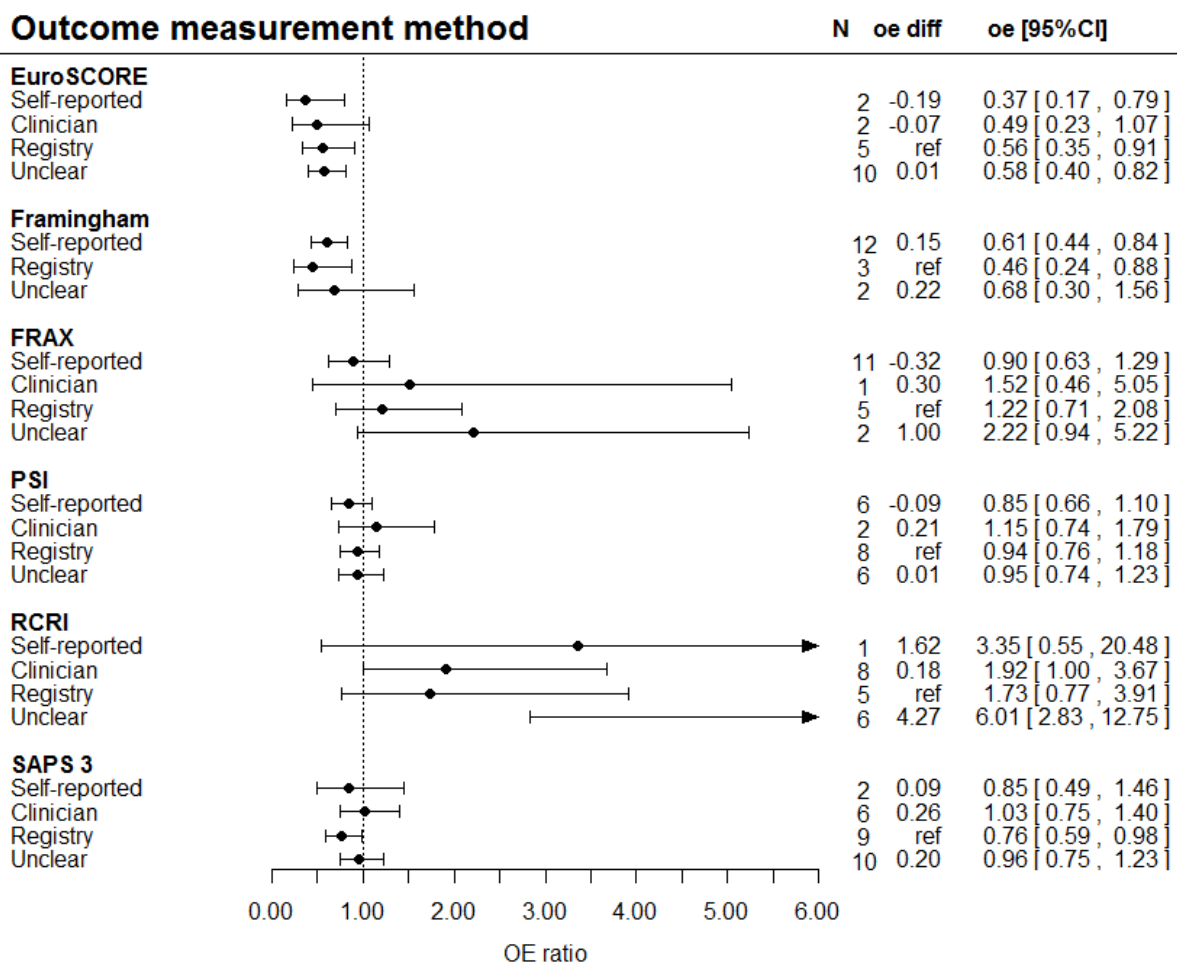


Damen et al. Empirical evidence on the impact of study characteristics and the performance of prediction models: a meta-epidemiological study

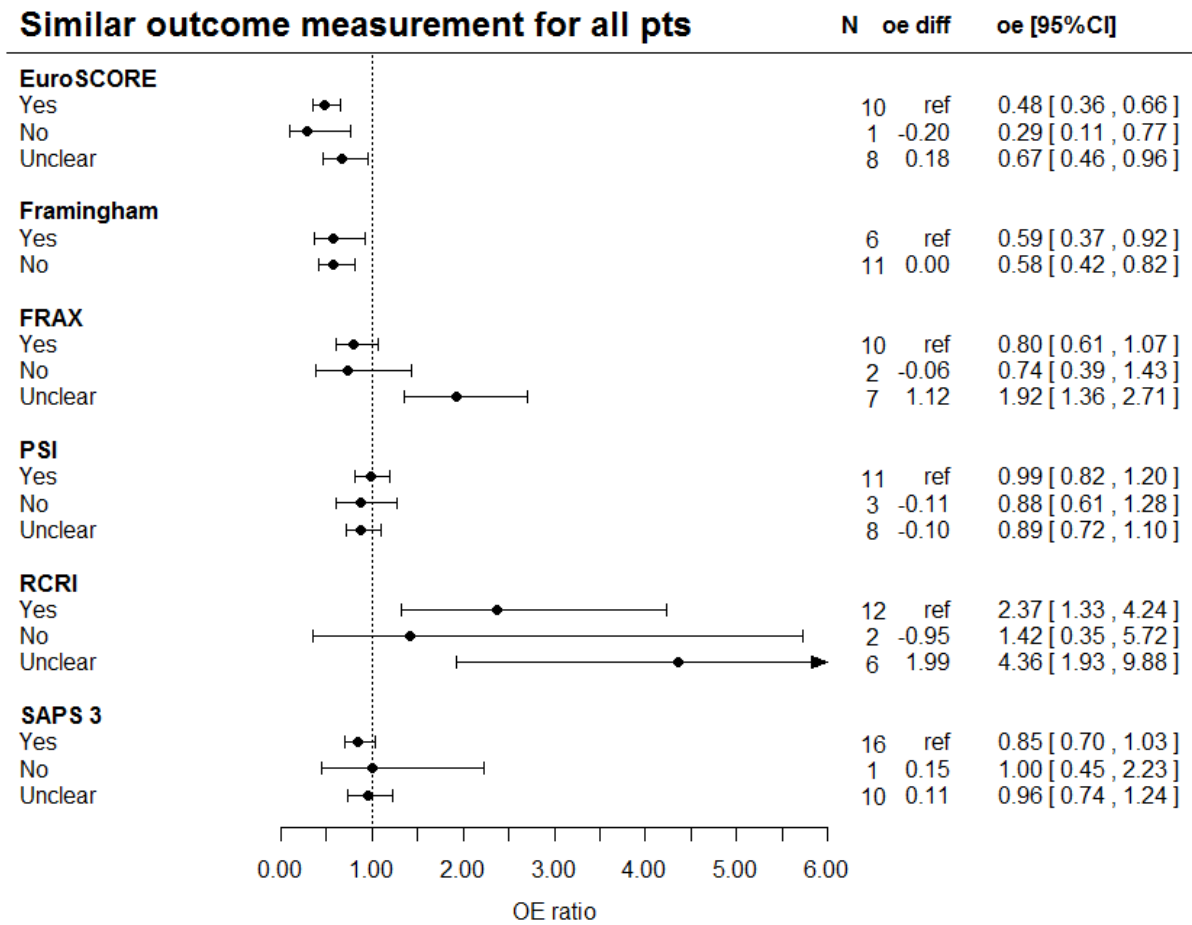


Damen et al. Empirical evidence on the impact of study characteristics and the performance of prediction models: a meta-epidemiological study

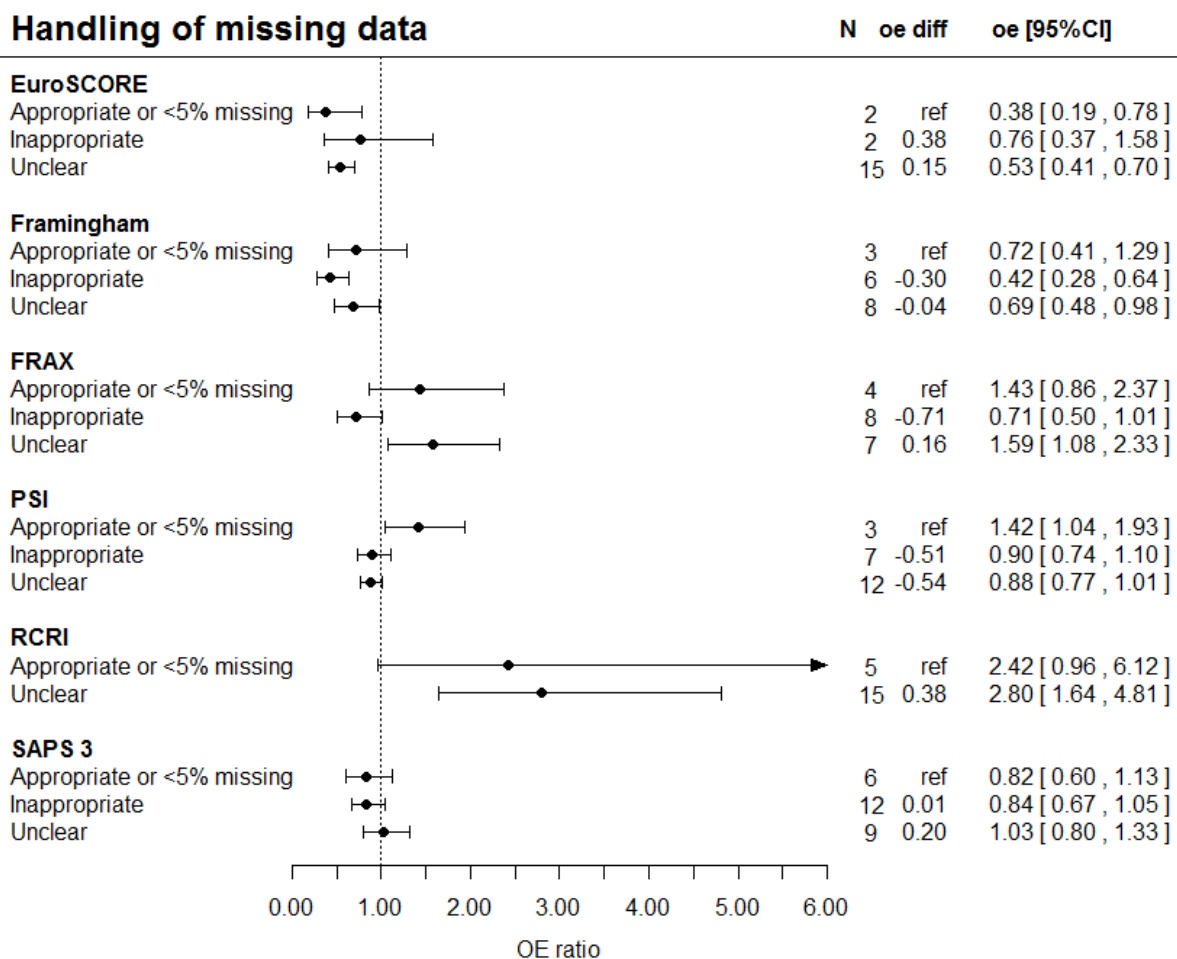
Outcome measurement method



Damen et al. Empirical evidence on the impact of study characteristics and the performance of prediction models: a meta-epidemiological study



Damen et al. Empirical evidence on the impact of study characteristics and the performance of prediction models: a meta-epidemiological study



Total OE ratio for categories of study characteristics, pooled using univariable meta-regression analyses per systematic review. N represents the number of external validation studies in a specific category. OE diff represents the difference in OE ratio with regard to a reference category (indicated with 'ref'). Dev: development, val: validation, incr: incremental value, pts: patients.

Damen et al. Empirical evidence on the impact of study characteristics and the performance of prediction models: a meta-epidemiological study

References

1. Snell KI, Ensor J, Debray TP, Moons KG, Riley RD. Meta-analysis of prediction model performance across multiple studies: Which scale helps ensure between-study normality for the C-statistic and calibration measures? *Stat Methods Med Res* 2017;962280217705678.
2. Newcombe RG. Confidence intervals for an effect size measure based on the Mann-Whitney statistic. Part 2: asymptotic methods and evaluation. *Stat Med* 2006;25(4):559-73.
3. Hanley JA, McNeil BJ. The meaning and use of the area under a receiver operating characteristic (ROC) curve. *Radiology* 1982;143(1):29-36.
4. Debray TP, Damen JA, Snell KI, Ensor J, Hooft L, Reitsma JB, et al. A guide to systematic review and meta-analysis of prediction model performance. *BMJ* 2017;356:i6460.
5. IntHout J, Ioannidis JP, Borm GF. The Hartung-Knapp-Sidik-Jonkman method for random effects meta-analysis is straightforward and considerably outperforms the standard DerSimonian-Laird method. *BMC Med Res Methodol* 2014;14:25.
6. Snell KI, Hua H, Debray TP, Ensor J, Look MP, Moons KG, et al. Multivariate meta-analysis of individual participant data helped externally validate the performance and implementation of a prediction model. *J Clin Epidemiol* 2015.
7. R: A language and environment for statistical computing [program]. Vienna, Austria: R Foundation for Statistical Computing, 2016.
8. Viechtbauer W. Conducting meta-analyses in R with the metafor package. *J Stat Softw* 2010;36(3):1-48.
9. Gasparrini A, Armstrong B, Kenward MG. Multivariate meta-analysis for non-linear and other multi-parameter associations. *Stat Med* 2012;31(29):3821-39.
10. Debray TP. Metamisc: Diagnostic and Prognostic Meta-Analysis. 2017.
11. Bates D, Mächler M, Bolker B, Walker S. Fitting linear mixed-effects models using lme4. *J Stat Softw* 2015;67(1).

Section/topic	Proposed item to be used in methodology research	Reported on page
TITLE		
Title	Identify the report as a meta-epidemiologic study.	1
ABSTRACT		
Structured summary	Provide a structured summary that includes the background of the topic, goal of the study, data sources, method of data selection, appraisal and synthesis methods, results, limitations, conclusions and implications of key findings.	3
INTRODUCTION		
Rationale	Describe the rationale for the meta-epidemiological study in the context of what is already known.	5
Objectives	Provide an explicit statement of the goal of the meta-epidemiological study and the hypothesis being empirically tested.	5
METHODS		
Protocol	Indicate if a protocol exists, if and where it can be accessed (eg, Web address). Registration of a protocol is not mandatory.	Available on request
Eligibility criteria	Specify study characteristics used as criteria for eligibility with a rationale.	7
Information sources	Describe all information sources (eg, databases with dates of coverage, contact with experts to identify additional studies, Internet searches) and search date.	7
Search	Present full electronic search strategy for at least one database, including any limits used, such that it could be repeated. Search is commonly not driven by a clinical question.	Supplement page 2
Study selection	Describe the process for selecting studies for inclusion (ie, how many reviewers selected studies, reviewing in duplicate or by single individuals).	7,8
Data collection process	Describe method of data extraction from reports (eg, piloted forms, independently, in duplicate) and any processes used for manipulating data or obtaining and confirming data from investigators.	8,9, Supplement page 3-6
Data items	List and define all variables for which data were sought and any assumptions and imputations made.	Supplement page 3-6
Risk of bias in individual studies	If risk of bias assessment of individual studies was relevant to the analysis, describe the items used and how this information is to be used during data synthesis.	Not assessed
Summary measures	State the principal summary measures (eg, ratio of risk ratios, difference in means) and explain its meaning and direction to readers.	9, Supplement page 7, 8
Synthesis of results	Describe the statistical or descriptive methods of synthesis including measures of consistency if relevant. If applicable, describe the development of statistical or simulation modelling based on theoretical background. Describe and justify assumptions and	9, Supplement page 7, 8

From: Murad MH, Wang Z. Guidelines for reporting meta-epidemiological methodology research. *Evid Based Med* 2017;22(4):139-42.

Section/topic	Proposed item to be used in methodology research	Reported on page
	computational approximations. Describe methods of additional analyses (eg, sensitivity or subgroup analyses, meta-regression), if done, indicating which were prespecified.	
RESULTS		
Study selection	Give numbers of studies assessed for eligibility and included in the study, with reasons for exclusions at each stage, ideally with a flow diagram. Present a measure of inter-reviewer agreement (eg, kappa statistic).	10, Figure 1
Study characteristics	For each study, present characteristics for which data were extracted and provide the citations. Clinical characteristics may not always be relevant.	Supplement page 9-12
Risk of bias within studies	If risk of bias assessment of individual studies was used in the meta-epidemiological analysis, report risk of bias indicators of each study to allow replication of findings.	Not assessed
Results of individual studies	Present data elements used in the meta-epidemiological analysis from each study (results of clinical outcomes may not be relevant).	Not done
Synthesis of results	Present results of statistical analysis done, including measures of precision and measures of consistency. Present validity of assumptions and fit of statistical or simulation modelling, if applicable.	11, 12, Figure 2-5, Supplement page 13-49
Additional analysis	Give results of additional analyses, if done (eg, sensitivity or subgroup analyses, meta-regression).	Not done
DISCUSSION		
Summary of evidence	Summarise the main findings and compare them with existing knowledge about the topic. The quality of evidence may not be relevant; however, investigators should describe their certainty in the results to readers.	13,14
Limitations	Discuss limitations at research methodology level (eg, likelihood of reporting or publication bias).	13,14
Conclusions	Provide general interpretation of the results and implications for future research. Provide any plausible impact on clinical practice.	16
FUNDING		
Funding	Describe sources of funding for the methodology research and role of funders.	17

From: Murad MH, Wang Z. Guidelines for reporting meta-epidemiological methodology research. Evid Based Med 2017;22(4):139-42.