

BMJ Open

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

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

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

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

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

BMJ Open

Curating a Knowledge Base for Individuals with Coinfection of HIV and SARS-CoV-2: A Study Protocol of EHR-based Data Mining and Clinical Implementation

Journal:	<i>BMJ Open</i>
Manuscript ID	bmjopen-2022-067204
Article Type:	Protocol
Date Submitted by the Author:	05-Aug-2022
Complete List of Authors:	Liang, Chen; University of South Carolina, Department of Health Services Policy and Management; University of South Carolina Weissman, Sharon; University of South Carolina, Olatosi, Bankole; University of South Carolina Arnold School of Public Health, Health Services, Policy and Management Poon, Eric ; Duke University Yarrington, Michael; Duke University, Li, Xiaoming; University of South Carolina Arnold School of Public Health
Keywords:	COVID-19, HIV & AIDS < INFECTIOUS DISEASES, Health informatics < BIOTECHNOLOGY & BIOINFORMATICS

SCHOLARONE™
Manuscripts

1
2
3 **Curating a Knowledge Base for Individuals with Coinfection of HIV and SARS-CoV-2: A Study**
4 **Protocol of EHR-based Data Mining and Clinical Implementation**
5
6

7 **Chen Liang^{1,2}, Sharon Weissman^{2,3}, Bankole Olatosi^{1,2}, Eric Poon⁴, Michael Yarrington⁴,**
8 **Xiaoming Li^{2,5}**
9

10
11 ¹ Department of Health Services Policy and Management, Arnold School of Public Health,
12 University of South Carolina

13 ² Big Data Health Science Center, University of South Carolina

14 ³ Department of Internal Medicine, School of Medicine, University of South Carolina

15 ⁴ Department of Medicine, School of Medicine, Duke University

16 ⁵ Department of Health Promotion Education and Behaviors, Arnold School of Public Health,
17 University of South Carolina
18
19

20 Correspondence:

21 Chen Liang, PhD

22 915 Greene St, Suite 347

23 Columbia, SC

24 cliang@mailbox.sc.edu

25 803 777 1836
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

Introduction: Despite a higher risk of severe COVID-19 disease in individuals with HIV, the interactions between SARS-CoV-2 and HIV infections remain unclear. To delineate these interactions, multi-center Electronic Health Records (EHR) hold existing promise to provide full-spectrum and longitudinal clinical data, demographics, and socio-behavioral data at individual-level. Presently, an EHR-based cohort for the HIV/SARS-CoV-2 coinfection has not been established; EHR integration and data mining methods tailored for studying the coinfection are urgently needed yet remain underdeveloped.

Methods and analysis: The overarching goal of this exploratory/developmental study is to establish an EHR-based cohort for individuals with HIV/SARS-CoV-2 coinfection and perform large-scale EHR-based data mining to examine the interactions between HIV and SARS-CoV-2 infections and systematically identify and validate factors contributing to the severe clinical course of the coinfection. We will use a nationwide EHR database in the United States, namely, National COVID Cohort Collaborative (N3C). Ultimately, collected clinical evidence will be implemented and used to pilot test a Clinical Decision Support (CDS) prototype to assist providers in screening and referral of at-risk patients in real-world clinics.

Ethics and dissemination: The study was approved by the institutional review boards at the University of South Carolina (Pro00121828) as non-human subject study. Study findings will be presented at academic conferences and published in peer-reviewed journals. This study will disseminate urgently needed clinical evidence for guiding clinical practice for individuals with the coinfection at Prisma Health, a healthcare system in collaboration.

Keywords

Acquired Immunodeficiency Syndrome, COVID-19, Electronic Health Records, Data Mining, Decision Support Systems

Strengths and limitations of this study

- This study will be among the first that systematically integrates HIV viral suppression status, antiretroviral therapy (ART) adherence, vaccination, socio-behavioral, and social determinants of health (SDOH) with full-spectrum clinical characteristics for individuals with the coinfection.
- Our methods can explain the role of temporal dependency among patients' underlying conditions, comorbidities, ART adherence, vaccine exposure, and received therapeutics in individuals' heterogeneous responses to the coinfection.
- Our methods support in-time prediction of coinfecting individuals' clinical outcomes, disease progression, and prognosis, and their risk factors.
- The proposed method is highly innovative in that it is designed to extract temporal sequences and temporal properties of every clinical event from EHR and is fully capable of embedding the temporal data into machine learning models.
- The study does not include comprehensive validity and usability testing of the proposed CDS prototype due to limited time for an exploratory/developmental study.

Introduction

The COVID-19 pandemic has cast a heavy burden on individuals with HIV infection. Based on data from 15,522 hospitalized patients with the coinfection of HIV and SARS-CoV-2 from 24 countries, a recent World Health Organization (WHO) report for the first time confirmed that HIV is a key risk factor for severe COVID-19¹. The severity of COVID-19 in individuals with HIV is correlated with certain comorbidities (e.g., type 2 diabetes mellitus, cardiovascular diseases, obesity, chronic obstructive pulmonary diseases, chronic kidney diseases, and some cancers) in which some comorbidities are more prevalent in people living with HIV (PLWH). Individuals with low CD4⁺ T-cell count (e.g., <200 cells/ μ L² or <500³ cells/ μ L) and unsuppressed viral load, and prolonged antiretroviral therapy (ART) exposure are associated with severe clinical course. These clinical facts are further complicated by the disrupted HIV healthcare services (e.g., access to HIV testing, ART, and distribution of PrEP and PEP).⁴

Despite a generally high risk of severe COVID-19 clinical course in PLWH, the interactions between SARS-CoV-2 and HIV infections remain unclear. First, several contradictory findings suggested the predominant role of comorbidities in severity of COVID-19 regardless of HIV infection⁵⁻⁸. Second, risk factors for the severe clinical course of the coinfection are undetermined because individuals with the same or similar severity level of COVID-19 show different clinical characteristics.⁴ Third, the role of ART adherence and HIV viral suppression status in the context of COVID-19 exposure is undetermined. These unsolved problems are attributed by several data and methodological gaps. For example, most existing studies are based on small-sample and single-center cohorts. Temporal sequences and patterns of clinical events (e.g., underlying conditions, comorbidities, diagnoses, ART-related visits and treatments) are understudied, which diminish the opportunities for understanding the etiology of multi-faceted HIV-associated comorbidities, their natural history, and their interactions with the current coinfection. Critical data components such as adherence to HIV treatment, viral suppression, social determinants of health (SDOH), COVID-19 vaccination, and socio-behavioral patterns (e.g., substance use/dependence) are closely related to disparities in HIV and SARS-CoV-2 infections but are understudied in part due to challenges in Electronic Health Records (EHR) data integration and phenotyping. The EHR hold existing promise to provide full-spectrum and longitudinal clinical data, demographics, and socio-behavioral data at the individual-level. However, we currently do not have an EHR-based cohort for individuals with HIV/SARS-CoV-2 coinfection; EHR integration and data mining tailored for studying the coinfection are urgently needed but not yet developed.

The overarching goal of this exploratory/developmental study is to establish an EHR-based cohort for individuals with HIV/SARS-CoV-2 coinfection and perform large-scale EHR-based data mining to examine the interactions between HIV and SARS-CoV-2 infections and systematically identify and validate factors contributing to the severe clinical course of the coinfection. Ultimately, collected clinical evidence will be implemented and used to pilot test a Clinical Decision Support (CDS) prototype to assist providers in screening and referral of at-risk patients in real-world clinics. We will approach this goal by pursuing the following tasks. First, we will extract comprehensive phenotypic traits (i.e., clinical characteristics, demographics, socio-behavioral patterns) and their temporal series and patterns from structured and unstructured EHR – National COVID Cohort Collaborative (N3C)⁹. To extract and model temporal series and

patterns of phenotypic traits, we will incorporate biomedical ontologies to develop a graphical model of EHR. Second, we will examine patterns and sequences of phenotypic traits for their predictive ability in clinical outcomes and disease prognosis. Major phenotypic traits to be examined include demographics, underlying conditions, comorbidities, CD4⁺ counts, viral suppression, ART procedures and medications, laboratory results for immune components and viral presence, treatments (e.g., procedures and medications), SDOH, and socio-behavioral patterns. We will develop machine-learning models to explore real-time predictive associations between these phenotypic traits and poor clinical outcomes and prognosis, including outcomes of the acute phase and post-acute sequelae of SARS-CoV-2 infection (PASC)¹⁰. Third, we will develop and pilot test a CDS prototype that delivers collected clinical evidence to providers through the Epic EHR system at Prisma Health. Predictive associations generated from the second task will be presented for providers to assist in screening patients at high risk of severe COVID-19 course. Outcomes to be measured include 1) the rate of identification and referral of individuals at high risk, 2) the rate of successful referral and clinical actions, and 3) system usability. The proposed study protocol will result in 1) a comprehensive knowledge base that details risk factors of severe clinical outcomes and disease prognosis in individuals with HIV/SARS-CoV-2 coinfection, 2) a prototype CDS that can identify patients at high risk and provide actionable clinical decisions. This work will provide time-sensitive public health implications: Clinical evidence for interactions between HIV and SARS-CoV-2 infections is desperately needed. This proposed EHR-based data mining offers a rapid and empirically grounded approach to collecting such evidence and to informing the design of prospective clinical trials that can focus on inflammatory pathways, biophysiological evidence of the coinfection, and socio-behavioral determinants.

Methods and Analysis

Data Description

We will use EHR from National COVID Cohort Collaborative (N3C). As of 3/2022, N3C has aggregated 12.4 million patients (4.5 million COVID-19 patients) from 50 states¹¹. The EHR are normalized by the Observational Health Data Sciences and Informatics (OHDSI)'s Observational Medical Outcomes Partnership (OMOP) Common Data Model (CDM)¹². EHR variables (individual level) span demographics, encounters, medical/social history, diagnoses, procedures, medication prescriptions, medication exposure (e.g., vaccines), laboratory tests/results, etc. Clinical notes have already been annotated in the CDM¹³. N3C has epidemiological and community data (population level) including temporal COVID-19 burden, vaccination, viral variance, health systems data, and geospatial data. N3C has pre-pandemic data since 2018 and peri-pandemic data to the present. All data were de-identified and updated every 2 weeks. We have the highest level of data access, which allows using patients' residential ZIP codes and dates of clinical events. As of August of 2021, N3C has at least 13 thousand adults have been diagnosed and/or have had laboratory-confirmed HIV. A pilot study shows that patients with coinfection have a high risk of hospitalization and mortality as compared to COVID-19 patients without HIV infection¹⁴.

EHR Data Modeling

We will remodel the EHR data extracted from N3C. Because Individuals' longitudinal health records are stored at distributed locations in EHR. Many clinical events do not have explicit and/or complete temporal information. Raw EHR data as such are of little value for understanding questions such as why individuals with certain coinfection present different clinical outcomes and disease prognosis.

We will locate relevant phenotypic traits from N3C by curating OMOP CDM concept sets, a procedure called electronic phenotyping. Using these concept sets, we will then retrieve and integrate phenotypic traits from N3C. At last, we will retrieve temporal information for clinical events and establish a graphical model¹⁵ to represent clinical events and their temporal information.

EHR phenotyping. Phenotyping is the process of identifying cohorts and variables from raw EHR¹⁶. Because we use OMOP CDM-based EHR, phenotyping is the process of finding "OMOP concepts" that correspond to specific cohorts (e.g., all patients with ART-related visits) or variables (e.g., myocardial infarction)¹⁷. An OMOP concept is a unique identifier that is mapped to diverse medical codes that have the same semantic meaning but maybe from different coding systems and EHR systems.

We will use standard phenotyping procedures. The logical procedures are as follows: 1) To identify existing OMOP concept sets (available in OHDSI's Atlas system and N3C) that can be used with minor revision. 2) If no appropriate concept sets available, we will follow the standardized phenotyping procedures^{16,18} and curate OMOP concepts from the Athena vocabulary repository, which allows retrieval of OMOP concepts¹⁹. 3) To validate the revised or newly developed concept sets by using EHR chart review that is performed independently by two clinical experts¹⁶.

Data retrieval and integration. To link external epidemiological and community-level data with individual-level EHR data, we will use individual patients' residential ZIP/county as the reference. We will use data imputation algorithms to impute missing values because population-to-individual-level data integration and missing values widely exist in EHR. For missing values in EHR, we will infer missing values based on semantic relationships of OMOP concepts. For geographical locations, we will use cities, counties, and states to infer locations at appropriate levels. For the rest of the missing values, we will use multiple imputation methods. Specifically, we will use selection models or pattern-mixture models for systematic missingness. We will also selectively use mean/median imputation, principal component analysis (PCA), singular value decomposition (SVD), k-nearest neighbor (kNN), least squares, expectation-maximization, and random forest.^{20,21} Data retrieval and integration will be implemented using SQL, R, and PySpark.

Graphical EHR model. We will develop a customized Graphical model to represent the temporal relations among patients' demographic, clinical, and socio-behavioral data. A Graphical model encodes clinical events as nodes and their semantic relations (including temporal relations) as edges. **Figure 1** shows the proposed design of the model. General modeling procedures include: First, we will retrieve time stamps of clinical events. Some clinical events have explicit time stamps as captured by structured EHR. Many others do not (e.g., patient-reported symptoms as documented in clinical notes, trimester and/or gestational age). For those that do not have explicit time stamps, we will infer the time of the event based on neighboring EHR data. For example, symptom onsets could be found in clinical notes (e.g., admission, discharge); trimester

1
2
3 and gestational age could be estimated by gestation-related diagnoses (e.g., Z3A: weeks of
4 gestation), procedures (e.g., ultrasound procedures), and clinical notes (e.g., last menstrual
5 period). Second, we will represent temporal information of clinical events by modeling the
6 occurrences of an event and the instantaneous impact of the event. The occurrences of an
7 event are resulted from the first step. The instantaneous impact of an event is formulated using
8 exponential kernel functions and association rules based on clinical observation. Intuitively, two
9 clinical events with a long interval in between would have less impact on one another, but this
10 can be overwritten by events that hold unique clinical meanings. Third, we will create recurrent
11 states for every clinical event by embedding the event, its semantic relations (i.e., edge in a
12 Graph) including the instantaneous impact of an event, and time. These recurrent states will
13 form a recurrent layer to be used for training machine learning models, which will be discussed
14 later.
15
16
17
18

19 **Machine Learning Modeling**

20 We will use supervised machine learning to examine patterns and sequences of phenotypic
21 traits for their predictive ability in clinical outcomes and disease prognosis. Existing studies
22 conclude differently on clinical outcomes among individuals with coinfection as well as the
23 factors correlated with these clinical outcomes. We provide two hypotheses. Hypothesis 1:
24 Individual patients respond differently to the coinfection. Hypothesis 2: Patients' clinical
25 outcomes and disease prognosis are attributed by the temporal dynamics of clinical events. If
26 these hypotheses are successfully tested, we will be able to delineate the impact of coinfection
27 on individuals' clinical outcomes. Therefore, we will customize recurrent neural network (RNN)
28 models to be used for predicting clinical outcomes and disease prognosis in real-time by
29 learning about patients' retrospective EHR at the individual level (personalized) as time
30 progresses. We adopt RNN for its unique advantage in capturing temporal dependencies of
31 data²². Trained RNN models will be tested for their performance where the best-performed
32 model will be identified for identifying patterns/sequences of phenotypic traits predictive of
33 clinical outcomes and prognosis.
34
35
36

37 **Cohort.** Based on the estimated >13 thousand patients with the coinfection in our dataset, we
38 will blend in controls (COVID-19 patients without HIV) for each output variable using a match
39 ratio of 1:2, stratified by sex, race/ethnicity, and age. For a possible occasions of small sample,
40 for example, death cases (n<1000 with coinfection), the alternative strategy is to create
41 synthetic cases to impute and balance the sample.
42

43 **Machine learning input.** A complete and longitudinal health history together with linked
44 external epidemiological and community-level data will be included as the input of machine
45 learning models by which the models can learn from the input to predict individuals' in-time
46 clinical outcomes and disease prognosis. We will include but are not limited to the following
47 phenotypic traits from N3C: demographics, SDOH, diagnoses, underlying conditions, vitals,
48 laboratory tests, procedures, medication prescriptions/dispensing, medication exposure (e.g.,
49 vaccine), and annotated clinical notes. Because there is no gold standard measure for ART
50 adherence, we will use "multiple measures" to estimate levels of ART adherence²³. Multiple
51 measures include medication events (inpatient dispensing), HIV-1 RNA copies (laboratory
52 results), and medication adherence data from clinical notes. For those with complete
53 medication adherence data we use the proportion of days covered (PDC) to categorize ART
54
55
56
57
58
59
60

1
2
3 adherence levels (e.g., < 50, 50-80, 80-85, 85-90, ≥90%)²⁴. We will categorize antiviral
4 medications (ARV) into Integrase inhibitor (INSTI)-based, non-nucleoside reverse transcriptase
5 inhibitor (NNRTI)-based, protease inhibitor (PI)-based, and other regimens. The approach to
6 measuring ART adherence using EHR has limitations, but the limitations can be mitigated by the
7 well-presented and large-scale national data. We collect both CD4⁺ counts as an indicator of
8 existing damage and plasma HIV-1 RNA copies as an indicator of projected disease progression.
9
10 ***Machine learning output.*** Clinical outcome measures as machine learning output include
11 inpatient admissions, length of stay (LOS), ICU admission, ICU LOS, comorbidities, and primary
12 discharge diagnosis²⁵. In addition to the measures within the acute phase of COVID-19, we will
13 also use symptoms, diagnoses, comorbidities, and PASC-associated-readmissions as outcome
14 and prognosis measures for individuals in the post-acute phase.

15
16 ***Model design.*** We will use RNN as the machine learning architecture to learn from patients'
17 longitudinal EHR and make the prediction of current and future clinical outcomes and disease
18 prognosis. We adopt RNN models because this neural network architecture is specialized for
19 capturing temporal dependency among event sequences. The RNN models will be trained to
20 learn from model input and to make the prediction of model output. With respect to the
21 embedding, we will include standard long short-term memory (LSTM) as well as phased LSTM
22 and other variants wherever appropriate²². We will use the bag-of-pattern matrix as the
23 baseline embedding method to be compared against LSTM, in which this baseline method does
24 not fully consider temporal dependency. To test against RNN, we will use Support Vector
25 Machine (SVM) as the benchmark algorithm, in which SVM is a well-performed kernel-based
26 algorithm²⁶ but does not take full advantage of temporal dependency (hypothesis 2) and
27 personalized health records (hypothesis 1). Because the nature of our machine learning output
28 is binary variables, the proposed machine learning tasks are essentially binary classification
29 tasks. We will use Python for machine learning modeling.

30
31 ***Model evaluation (internal validity).*** To test the effectiveness of the prediction model, we will
32 use 10-fold cross validation. With respect to evaluation metrics, we will use the F score,
33 precision, recall, and the Area Under the Receiver Operating Characteristic (AUC) to assess the
34 models' predictive performance. We expect the F score, assuming balanced data, to reach a
35 minimum of 0.8. If trained models fail to meet this expectation, alternative strategies include
36 manually adding features handpicked by researchers after error analysis of models.

41 42 ***Clinical Decision Support System***

43 Based on the automatic clinical outcomes and disease prognosis prediction model, we will
44 design and implement a clinical decision support (CDS) prototype in collaboration with Prisma
45 Health clinics. The proposed CDS prototype is anticipated to assist providers in screening and
46 identifying patients who are at high risk of worse COVID-19 clinical outcomes (see **Table 1**), and
47 worse disease prognosis [including individuals with post-acute sequelae of SARS-CoV-2
48 infection (PASC)]. Specifically, the CDS will identify individuals with a high risk of disease
49 progression from their current clinical state (e.g., not hospitalized, hospitalized, post-acute
50 phase) by learning from the trained machine-learning models. The effectiveness of the CDS
51 demonstrates the external validity of the internally validated predictive model and will be
52 assessed by 1) appropriate identification for at-risk individuals, 2) appropriate clinical actions,
53 and 3) CDS system usability.

Table 1. Patient state according to WHO clinical progression scale.

Patient state	Clinical characteristics	Severity score
Uninfected	Uninfected, no viral RNA detected	0
Ambulatory mild disease	Asymptomatic, viral RNA detected	1
	Symptomatic, independent	2
	Symptomatic, assistance needed	3
Hospitalized, moderate disease	Hospitalized, no oxygen therapy	4
	Hospitalized, oxygen by mask or nasal prongs	5
Hospitalized, severe disease	Hospitalized, oxygen by NIV or high flow	6
	Intubation and mechanical ventilation, $pO_2/FiO_2 \geq 150$ or $SpO_2/FiO_2 \geq 200$	7
	Mechanical ventilation, $pO_2/FiO_2 < 150$ ($SpO_2/FiO_2 < 200$) or vasopressors	8
	Mechanical ventilation, $pO_2/FiO_2 < 150$ and vasopressors, dialysis, or ECMO	9
Dead	Dead	10

CDS workflow. The proposed CDS is a hybrid of knowledge-based and non-knowledge-based system²⁷. It has 1) a machine-learning-based prediction module (non-knowledge-based) for identifying high-risk patients and 2) a provider-curated medical logic module (knowledge-based) for generating clinical actions for identified high-risk patients. The CDS testing takes place in a retrospective way (i.e., using retrospective EHR).

Cohort definition and data collection. Using retrospective EHR data (2-year baseline~2023) from Prisma Health's Epic system, we will first group the existing PLWH who have COVID-19 (sampling $n > 500$) based on their state at the point of CDS screening. The patient states include 1) ambulatory patients with COVID-19, 2) hospitalized patients for COVID-19 with moderate disease, 3) hospitalized patients for COVID-19 with severe disease, and 4) post-acute phase of COVID-19 (i.e., from beyond 4 weeks after symptom onset)²⁸. See **Table 1** for definitions of states 1-3 based on WHO's clinical progression scale for COVID-19.

Prediction module. For patients in each state, we will use the trained machine-learning model to learn from previous medical records and predict worsening clinical outcomes as time progresses (i.e., acute, and post-acute phases every 3 months). The prediction will include primary COVID-19 clinical outcomes (**Table 2**) developed by the WHO working group on the Clinical Characterisation and Management of COVID-19²⁹.

Table 2. Key clinical outcome measures.

Organ dysfunction
• Murray score
• Sequential organ failure assessment score, multiple organ dysfunction score
• Acute coronary syndrome; arrhythmias
• Delirium

Comorbidities <ul style="list-style-type: none"> Pulmonary, cardiovascular, renal, neurological, etc.
Secondary infection <ul style="list-style-type: none"> Bacterial, viral
Biochemical parameters <ul style="list-style-type: none"> C-reactive protein, D-dimers, IL-6, and ferritin serum concentrations, and leucocyte counts
Radiological findings <ul style="list-style-type: none"> Chest CT scan, chest x-ray
Duration of intervention <ul style="list-style-type: none"> Inpatient admission, length of stay (LOS) ICU admission, ICU LOS Ventilation Organ support or hospital-free days
Pregnancy outcomes <ul style="list-style-type: none"> Preterm delivery, miscarriage Fetal status Severe maternal morbidity (SMM) measures
Mortality <ul style="list-style-type: none"> All-cause mortality at hospital discharge
Quality of life <ul style="list-style-type: none"> Longer term survival and primary diagnoses for readmission (post-acute phase)

Medical logic module. Patients identified by the CDS to have an increased risk of worse clinical outcomes will be reviewed and discussed by two providers who are specialized in HIV and COVID-19. First, the providers will generate gold-standard judgment on whether a patient is correctly identified by the prediction module, which will later be used for assessing the effectiveness of CDS. Second, the providers will generate appropriate clinical actions upon chart review. These clinical actions will be made up to date with the “NIH Guidance for COVID-19 and People with HIV”, including treatment options based on cohorts and risk factors, medication reconciliation considering ART regimens, consultation with specialists for multi-organ system complications and PASC, referrals, and outreach³⁰. Providers’ decision-making processes will be programmed using Arden Syntax (v3) or Clinical Quality Language³¹ in the knowledge base, which is determined by specific Epic EHR data model.

Effectiveness of CDS (external validity). There are two evaluation metrics: 1) Appropriate identification for individuals at high risk for adverse clinical outcomes (**Table 2**) by comparing model-identified cases against the gold standard generated from chart review. We will use F measure (>0.8), AUC, precision, and recall for assessment; 2) Appropriate clinical actions using a quasi-experimental design. We will compare outcomes of patients who naturally used the medical logic module-suggested care against those who did not (n=100 each). The outcomes include but are not limited to readmissions (e.g., same day, 7-, 14-, 30- days), healthcare utilization (e.g., LOS, ER/observation visits, ICU admission). We will use mixed-effect generalized regression models to estimate model effectiveness wherever appropriate.

1
2
3 *Usability testing.* We will assess CDS usability by adopting the “think aloud” protocol³². The two
4 providers from Prisma Health will participate in the test. Each one will be presented with
5 randomly selected EHR (n=5 at-risk cases + n=5 control cases) along with the CDS output. In
6 each case, participants will be instructed to verbalize their reasoning procedures (e.g.,
7 phenotypic traits from EHR that can be used in the reasoning, logic flow) towards identifying
8 the at-risk patients and corresponding clinical decisions. Sessions are audio-recorded and will
9 then be coded (e.g., by content, understandability, navigation, workflow, visibility, and
10 usability) independently by two researchers for downstream analyses.
11
12

13 14 **Patient and Public Involvement**

15 No patient involved.
16
17
18

19 **Ethics and Dissemination**

20 The study was approved by the institutional review boards at the University of South Carolina
21 (Pro00121828) as non-human subject study.
22

23 This study will result in a comprehensive knowledge base that documents clinical outcomes and
24 disease prognosis for individuals with the coinfection, their risk factors (e.g., underlying
25 conditions, ART adherence, comorbidities, socio-behavioral), and their responses to
26 therapeutics. This study will also result in a prototype CDS that can identify patients at high risk
27 of worsening clinical outcomes and prognosis in real-time. These results are generalizable and
28 will form a foundation for developing comprehensive real-world CDS systems for
29 implementation in state-wide and national HIV and COVID-19 clinics.
30

31 Study findings will be presented at academic conferences and published in peer-reviewed
32 journals. This study will disseminate urgently needed clinical evidence for guiding clinical
33 practice for individuals with the coinfection at Prisma Health.
34
35
36
37

38 **Authors' Contributions**

39 CL conceived the study design and drafted the manuscript. CL completed the preliminary data
40 analysis. SW, BO, EP, MY, and XL contributed critical edits to the manuscript. All authors
41 reviewed and approved the manuscript.
42
43
44

45 **Funding Statement**

46 Research reported in this publication was supported by the National Institute Of Allergy And
47 Infectious Diseases of the National Institutes of Health under Award Number R21AI170171. The
48 content is solely the responsibility of the authors and does not necessarily represent the official
49 views of the National Institutes of Health.
50
51
52

53 **Competing Interest Statement**

54 None declared.
55
56
57
58
59
60

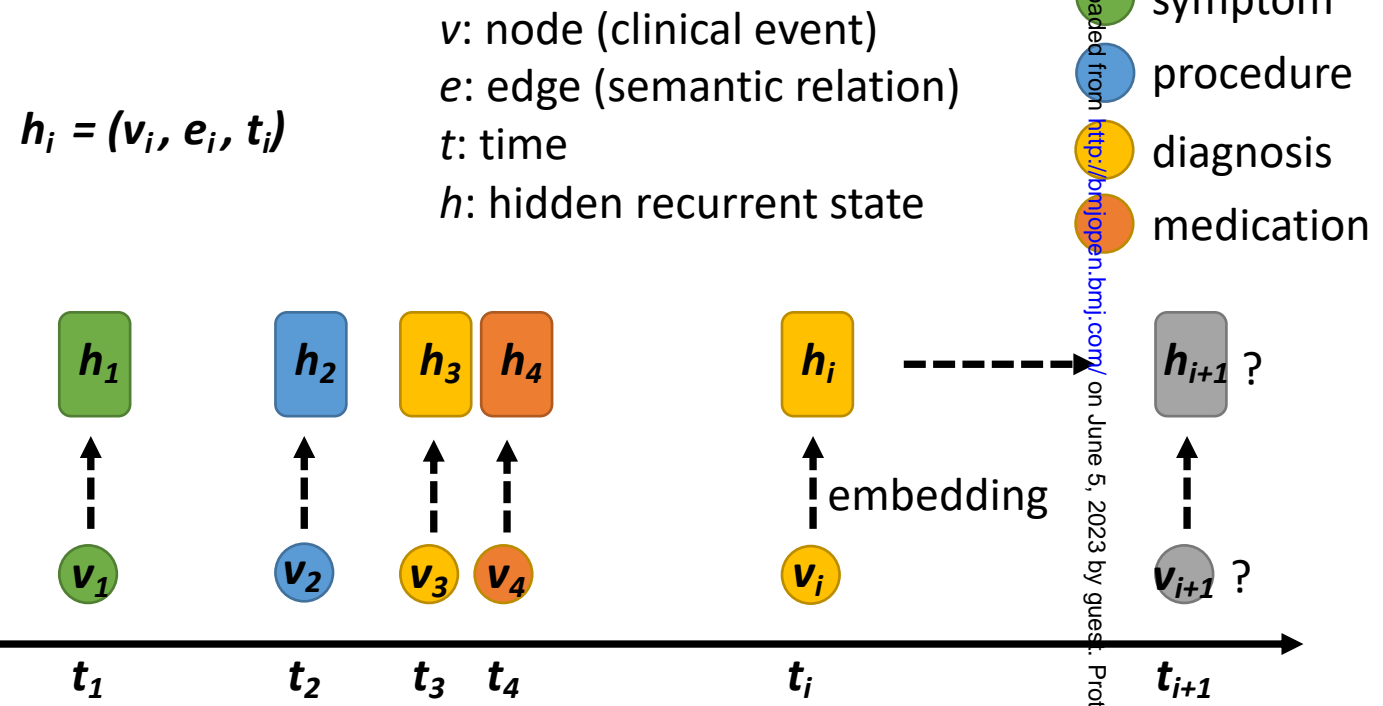
References

1. World Health Organization. *Clinical Features and Prognostic Factors of COVID-19 in People Living with HIV Hospitalized with Suspected or Confirmed SARS-CoV-2 Infection.*; 2021.
2. Dandachi D, Geiger G, Montgomery MW, et al. Characteristics, Comorbidities, and outcomes in a multicenter registry of patients with human immunodeficiency virus and coronavirus disease 2019. *Clin Infect Dis*. Published online 2020.
3. Braunstein SL, Lazar R, Wahnich A, Daskalakis DC, Blackstock OJ. COVID-19 infection among people with HIV in New York City: A population-level analysis of linked surveillance data. *Clinical Infectious Diseases: An Official Publication of the Infectious Diseases Society of America*. Published online 2020.
4. Eisinger RW, Lerner AM, Fauci AS. Human Immunodeficiency Virus/AIDS in the Era of Coronavirus Disease 2019: A Juxtaposition of 2 Pandemics. *The Journal of Infectious Diseases*. Published online 2021.
5. Cooper TJ, Woodward BL, Alom S, Harky A. Coronavirus disease 2019 (COVID-19) outcomes in HIV/AIDS patients: a systematic review. *HIV Med*. 2020;21(9):567-577.
6. Calza L, Bon I, Tadolini M, et al. COVID-19 in patients with HIV-1 infection: a single-centre experience in northern Italy. *Infection*. 2021;49(2):333-337.
7. Costenaro P, Minotti C, Barbieri E, Giaquinto C, Donà D. SARS-CoV-2 infection in people living with HIV: a systematic review. *Reviews in Medical Virology*. 2021;31(1):1-12.
8. Park LS, Rentsch CT, Sigel K, et al. COVID-19 in the largest US HIV cohort. *AIDS*. 2020;2020:23rd.
9. Haendel MA, Chute CG, Gersing K. The National COVID Cohort Collaborative (N3C): rationale, design, infrastructure, and deployment. *Journal of the American Medical Informatics Association*. Published online 2020.
10. Deer RR, Rock MA, Vasilevsky N, et al. Characterizing long COVID: deep phenotype of a complex condition. *EBioMedicine*. 2021;74:103722.
11. Datavent. COVID-19 Research Database. Accessed February 20, 2021. <https://covid19researchdatabase.org/>
12. OHDSI community. Observational Health Data Sciences and Informatics Common Data Model.
13. N3C. COVID-19 Clinical Data Warehouse Data Dictionary.
14. Yang X, Zhang J, Guo S, Olatosi B, Weissman S, Li X. The role of HIV infection in the clinical spectrum of COVID-19: a population-based cohort analysis based on US National COVID Cohort Collaborative (N3C) Enclave data. *Available at SSRN 3860395*. Published online 2021.
15. Liu C, Wang F, Hu J, Xiong H. Temporal phenotyping from longitudinal electronic health records: A graph based framework. In: *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining.* ; 2015:705-714.
16. Banda JM, Seneviratne M, Hernandez-Boussard T, Shah NH. Advances in electronic phenotyping: from rule-based definitions to machine learning models. *Annu Rev Biomed Data Sci*. 2018;1:53-68.
17. Richesson RL, Hammond WE, Nahm M, et al. Electronic health records based phenotyping in next-generation clinical trials: a perspective from the NIH Health Care

- 1
2
3 Systems Collaboratory. *Journal of the American Medical Informatics Association*.
4 2013;20(e2):e226-e231.
5
6 18. Weng C, Shah NH, Hripcsak G. Deep phenotyping: Embracing complexity and
7 temporality—Towards scalability, portability, and interoperability. *Journal of Biomedical*
8 *Informatics*. 2020;105:103433.
9
10 19. OHDSI Athena standard vocabularies. Accessed September 1, 2021.
11 <https://www.ohdsi.org/analytic-tools/athena-standardized-vocabularies/>
12
13 20. Luo Y. Evaluating the state of the art in missing data imputation for clinical data. *Briefings*
14 *in Bioinformatics*. 2022;23(1):bbab489.
15
16 21. Li J, Yan XS, Chaudhary D, et al. Imputation of missing values for electronic health record
17 laboratory data. *NPJ Digit Med*. 2021;4(1):1-14.
18
19 22. Goodfellow I, Bengio Y, Courville A, Bengio Y. *Deep Learning*. Vol 1. MIT press Cambridge;
20 2016.
21
22 23. Castillo-Mancilla JR, Haberer JE. Adherence measurements in HIV: new advancements in
23 pharmacologic methods and real-time monitoring. *Current Hiv/aids Reports*.
24 2018;15(1):49-59.
25
26 24. Byrd KK, Hou JG, Hazen R, et al. Antiretroviral adherence level necessary for HIV viral
27 suppression using real-world data. *J Acquir Immune Defic Syndr*. 2019;82(3):245.
28
29 25. Lavery AM, Preston LE, Ko JY, et al. Characteristics of Hospitalized COVID-19 Patients
30 Discharged and Experiencing Same-Hospital Readmission—United States, March--August
31 2020. *Morbidity and Mortality Weekly Report*. 2020;69(45):1695.
32
33 26. Murphy KP. *Machine Learning: A Probabilistic Perspective*. MIT press; 2012.
34
35 27. Shiffman RN, Wright A. Evidence-based clinical decision support. *Yearb Med Inform*.
36 2013;22(01):120-127.
37
38 28. Nalbandian A, Sehgal K, Gupta A, et al. Post-acute COVID-19 syndrome. *Nat Med*.
39 Published online 2021:1-15.
40
41 29. Marshall JC, Murthy S, Diaz J, et al. A minimal common outcome measure set for COVID-
42 19 clinical research. *The Lancet Infectious Diseases*. 2020;20(8):e192–e197.
43
44 30. Guidelines Working Groups of the NIH Office of AIDS Research Advisory Council.
45 Guidance for COVID-19 and People with HIV.
46
47 31. Hripcsak G, Clayton P, Pryor T, Haug P, Wigertz O, der Lei J. The Arden syntax for medical
48 logic modules. In: *Proceedings. Symposium on Computer Applications in Medical Care*. ;
49 1990:200-204.
50
51 32. Li AC, Kannry JL, Kushniruk A, et al. Integrating usability testing and think-aloud protocol
52 analysis with “near-live” clinical simulations in evaluating clinical decision support. *Int J*
53 *Med Inform*. 2012;81(11):761-772.
54
55
56
57
58
59
60

1
2
3 Figure 1. EHR model design.
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

067204 on 13 September 2022. Downloaded from <http://bmjopen.bmj.com/> on June 5, 2023 by guest. Protected by copyright.



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



CONSORT 2010 checklist of information to include when reporting a randomised trial*

Section/Topic	Item No	Checklist item	Reported on page No
Title and abstract			
	1a	Identification as a randomised trial in the title	NA
	1b	Structured summary of trial design, methods, results, and conclusions (for specific guidance see CONSORT for abstracts)	2
Introduction			
Background and objectives	2a	Scientific background and explanation of rationale	4-5
	2b	Specific objectives or hypotheses	4-5
Methods			
Trial design	3a	Description of trial design (such as parallel, factorial) including allocation ratio	NA
	3b	Important changes to methods after trial commencement (such as eligibility criteria), with reasons	NA
Participants	4a	Eligibility criteria for participants	5, 8
	4b	Settings and locations where the data were collected	5
Interventions	5	The interventions for each group with sufficient details to allow replication, including how and when they were actually administered	NA
Outcomes	6a	Completely defined pre-specified primary and secondary outcome measures, including how and when they were assessed	11
	6b	Any changes to trial outcomes after the trial commenced, with reasons	NA
Sample size	7a	How sample size was determined	8
	7b	When applicable, explanation of any interim analyses and stopping guidelines	NA
Randomisation:			
Sequence generation	8a	Method used to generate the random allocation sequence	NA
	8b	Type of randomisation; details of any restriction (such as blocking and block size)	NA
Allocation concealment mechanism	9	Mechanism used to implement the random allocation sequence (such as sequentially numbered containers), describing any steps taken to conceal the sequence until interventions were assigned	NA
Implementation	10	Who generated the random allocation sequence, who enrolled participants, and who assigned participants to interventions	NA
Blinding	11a	If done, who was blinded after assignment to interventions (for example, participants, care providers, those	NA

		assessing outcomes) and how	
	11b	If relevant, description of the similarity of interventions	NA
Statistical methods	12a	Statistical methods used to compare groups for primary and secondary outcomes	NA
	12b	Methods for additional analyses, such as subgroup analyses and adjusted analyses	NA
Results			
Participant flow (a diagram is strongly recommended)	13a	For each group, the numbers of participants who were randomly assigned, received intended treatment, and were analysed for the primary outcome	NA
	13b	For each group, losses and exclusions after randomisation, together with reasons	NA
Recruitment	14a	Dates defining the periods of recruitment and follow-up	NA
	14b	Why the trial ended or was stopped	NA
Baseline data	15	A table showing baseline demographic and clinical characteristics for each group	NA
Numbers analysed	16	For each group, number of participants (denominator) included in each analysis and whether the analysis was by original assigned groups	8
Outcomes and estimation	17a	For each primary and secondary outcome, results for each group, and the estimated effect size and its precision (such as 95% confidence interval)	NA
	17b	For binary outcomes, presentation of both absolute and relative effect sizes is recommended	NA
Ancillary analyses	18	Results of any other analyses performed, including subgroup analyses and adjusted analyses, distinguishing pre-specified from exploratory	11
Harms	19	All important harms or unintended effects in each group (for specific guidance see CONSORT for harms)	NA
Discussion			
Limitations	20	Trial limitations, addressing sources of potential bias, imprecision, and, if relevant, multiplicity of analyses	NA
Generalisability	21	Generalisability (external validity, applicability) of the trial findings	9-11
Interpretation	22	Interpretation consistent with results, balancing benefits and harms, and considering other relevant evidence	9-11
Other information			
Registration	23	Registration number and name of trial registry	NA
Protocol	24	Where the full trial protocol can be accessed, if available	NA
Funding	25	Sources of funding and other support (such as supply of drugs), role of funders	12

*We strongly recommend reading this statement in conjunction with the CONSORT 2010 Explanation and Elaboration for important clarifications on all the items. If relevant, we also recommend reading CONSORT extensions for cluster randomised trials, non-inferiority and equivalence trials, non-pharmacological treatments, herbal interventions, and pragmatic trials. Additional extensions are forthcoming: for those and for up to date references relevant to this checklist, see www.consort-statement.org.

BMJ Open

Curating a Knowledge Base for Individuals with Coinfection of HIV and SARS-CoV-2: A Study Protocol of EHR-based Data Mining and Clinical Implementation

Journal:	<i>BMJ Open</i>
Manuscript ID	bmjopen-2022-067204.R1
Article Type:	Protocol
Date Submitted by the Author:	24-Aug-2022
Complete List of Authors:	Liang, Chen; University of South Carolina, Department of Health Services Policy and Management; University of South Carolina Weissman, Sharon; University of South Carolina, Olatosi, Bankole; University of South Carolina Arnold School of Public Health, Health Services, Policy and Management Poon, Eric ; Duke University Yarrington, Michael; Duke University, Li, Xiaoming; University of South Carolina Arnold School of Public Health
Primary Subject Heading:	Infectious diseases
Secondary Subject Heading:	Health informatics, HIV/AIDS, Respiratory medicine
Keywords:	COVID-19, HIV & AIDS < INFECTIOUS DISEASES, Health informatics < BIOTECHNOLOGY & BIOINFORMATICS

SCHOLARONE™
Manuscripts

1
2
3 **Curating a Knowledge Base for Individuals with Coinfection of HIV and SARS-CoV-2: A Study**
4 **Protocol of EHR-based Data Mining and Clinical Implementation**
5
6

7 **Chen Liang^{1,2,3}, Sharon Weissman^{2,4}, Bankole Olatosi^{1,2,3}, Eric Poon⁵, Michael Yarrington⁵,**
8 **Xiaoming Li^{2,3,6}**
9

10
11 ¹ Department of Health Services Policy and Management, Arnold School of Public Health,
12 University of South Carolina

13 ² Big Data Health Science Center, University of South Carolina

14 ³ South Carolina SmartState Center for Healthcare Quality, University of South Carolina

15 ⁴ Department of Internal Medicine, School of Medicine, University of South Carolina

16 ⁵ Department of Medicine, School of Medicine, Duke University

17 ⁶ Department of Health Promotion Education and Behaviors, Arnold School of Public Health,
18 University of South Carolina
19
20

21 Correspondence:

22 Chen Liang, PhD

23 915 Greene St, Suite 347

24 Columbia, SC

25 cliang@mailbox.sc.edu

26 803 777 1836
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

Introduction: Despite a higher risk of severe COVID-19 disease in individuals with HIV, the interactions between SARS-CoV-2 and HIV infections remain unclear. To delineate these interactions, multi-center Electronic Health Records (EHR) hold existing promise to provide full-spectrum and longitudinal clinical data, demographics, and socio-behavioral data at individual-level. Presently, a comprehensive EHR-based cohort for the HIV/SARS-CoV-2 coinfection has not been established; EHR integration and data mining methods tailored for studying the coinfection are urgently needed yet remain underdeveloped.

Methods and analysis: The overarching goal of this exploratory/developmental study is to establish an EHR-based cohort for individuals with HIV/SARS-CoV-2 coinfection and perform large-scale EHR-based data mining to examine the interactions between HIV and SARS-CoV-2 infections and systematically identify and validate factors contributing to the severe clinical course of the coinfection. We will use a nationwide EHR database in the United States, namely, National COVID Cohort Collaborative (N3C). Ultimately, collected clinical evidence will be implemented and used to pilot test a Clinical Decision Support (CDS) prototype to assist providers in screening and referral of at-risk patients in real-world clinics.

Ethics and dissemination: The study was approved by the institutional review boards at the University of South Carolina (Pro00121828) as non-human subject study. Study findings will be presented at academic conferences and published in peer-reviewed journals. This study will disseminate urgently needed clinical evidence for guiding clinical practice for individuals with the coinfection at Prisma Health, a healthcare system in collaboration.

Keywords

Acquired Immunodeficiency Syndrome, COVID-19, Electronic Health Records, Data Mining, Decision Support Systems

Strengths and limitations of this study

- This study will be among the first that systematically integrates HIV viral suppression status, antiretroviral therapy (ART) adherence, vaccination, socio-behavioral, and social determinants of health (SDOH) with full-spectrum clinical characteristics for individuals with HIV/SARS-CoV-2 coinfection.
- Our methods can explain the role of temporal dependency among patients' underlying conditions, comorbidities, ART adherence, vaccine exposure, and received therapeutics in individuals' heterogeneous responses to the coinfection.
- Our methods support real-time prediction of coinfecting individuals' clinical outcomes, disease progression, and prognosis, and risk factors of adverse events.
- The proposed methods are highly innovative in that they are designed to extract temporal sequences and temporal properties of every clinical event from EHR and is fully capable of embedding the temporal data into machine learning models.
- The study does not include comprehensive validity and usability testing of the proposed Clinical Decision Support (CDS) prototype due to limited time for an exploratory/developmental study.

Introduction

The COVID-19 pandemic has cast a heavy burden on individuals with HIV infection. Based on data from 15,522 hospitalized patients with the coinfection of HIV and SARS-CoV-2 from 24 countries, a recent World Health Organization (WHO) report for the first time confirmed that HIV is a key risk factor for severe COVID-19¹. The severity of COVID-19 in individuals with HIV is correlated with certain comorbidities (e.g., type 2 diabetes mellitus, cardiovascular diseases, obesity, chronic obstructive pulmonary diseases, chronic kidney diseases, and some cancers) in which some comorbidities are more prevalent in people living with HIV (PLWH). Individuals with low CD4⁺ T-cell count (e.g., <200 cells/ μ L² or <500³ cells/ μ L) and unsuppressed viral load, and prolonged antiretroviral therapy (ART) exposure are associated with severe clinical course. These clinical facts are further complicated by the disrupted HIV healthcare services [e.g., access to HIV testing, ART, and distribution of pre-exposure prophylaxis (PrEP) and post-exposure prophylaxis (PEP)].⁴

Despite a generally high risk of severe COVID-19 clinical course in PLWH, the interactions between SARS-CoV-2 and HIV infections remain unclear. First, several contradictory findings suggested the predominant role of comorbidities in severity of COVID-19 regardless of HIV infection⁵⁻⁸. Second, risk factors for the severe clinical course of the coinfection are undetermined because individuals with the same or similar severity level of COVID-19 show different clinical characteristics.⁴ Third, the role of ART adherence and HIV viral suppression status in the context of COVID-19 exposure is undetermined. These unsolved problems are attributed by several data and methodological gaps. For example, most existing studies are based on small-sample and single-center cohorts. Temporal sequences and patterns of clinical events (e.g., underlying conditions, comorbidities, diagnoses, ART-related visits and treatments) are understudied, which diminish the opportunities for understanding the etiology of multi-faceted HIV-associated comorbidities, their natural history, and their interactions with the current coinfection. Critical data components such as adherence to HIV treatment, viral suppression, social determinants of health (SDOH), COVID-19 vaccination, and socio-behavioral patterns (e.g., substance use/dependence) are closely related to disparities in HIV and SARS-CoV-2 infections but are understudied in part due to the challenges in Electronic Health Records (EHR) data integration and phenotyping. EHR hold existing promise to provide full-spectrum and longitudinal clinical data, demographics, and socio-behavioral data at the individual-level. However, currently we do not have a comprehensive EHR-based cohort for individuals with HIV/SARS-CoV-2 coinfection; EHR integration and data mining tailored for studying the coinfection are urgently needed but are not yet developed.

The overarching goal of this exploratory/developmental study is to establish an EHR-based cohort for individuals with HIV/SARS-CoV-2 coinfection and perform large-scale EHR-based data mining to examine the interactions between HIV and SARS-CoV-2 infections and systematically identify and validate factors contributing to the severe clinical course of the coinfection. Ultimately, collected clinical evidence will be implemented and used to pilot test a Clinical Decision Support (CDS) prototype to assist providers in screening and referral of at-risk patients in real-world clinics. We will approach this goal by pursuing the following tasks. First, we will extract comprehensive phenotypic traits (i.e., clinical characteristics, demographics, socio-behavioral patterns) and their temporal series and patterns from structured and unstructured

1
2
3 EHR – National COVID Cohort Collaborative (N3C)⁹. To extract and model temporal series and
4 patterns of phenotypic traits, we will incorporate biomedical ontologies to develop a graphical
5 model of EHR. Second, we will examine patterns and sequences of phenotypic traits for their
6 predictive ability in clinical outcomes and disease prognosis. Major phenotypic traits to be
7 examined include demographics, underlying conditions, comorbidities, CD4⁺ counts, viral
8 suppression, ART procedures and medications, laboratory results for immune components and
9 viral presence, treatments (e.g., procedures and medications), SDOH, and socio-behavioral
10 patterns. We will develop machine-learning models to explore real-time predictive associations
11 between these phenotypic traits and poor clinical outcomes and prognosis, including outcomes
12 of the acute phase of COVID-19 and post-acute sequelae of SARS-CoV-2 infection (PASC)¹⁰.
13 Third, we will develop and pilot test a CDS prototype that delivers collected clinical evidence to
14 providers through the Epic EHR system at Prisma Health. Predictive associations generated
15 from the second task will be presented for providers to assist in screening patients at high risk
16 of severe COVID-19 course. Outcomes to be measured include 1) the rate of identification and
17 referral of individuals at high risk of poor clinical outcomes, 2) the rate of successful referral and
18 clinical actions, and 3) system usability. The proposed study protocol will result in 1) a
19 comprehensive knowledge base that details risk factors of severe clinical outcomes and disease
20 prognosis in individuals with HIV/SARS-CoV-2 coinfection, 2) a prototype CDS that can identify
21 patients at high risk and provide actionable clinical decisions. This work will provide time-
22 sensitive public health implications: Clinical evidence for interactions between HIV and SARS-
23 CoV-2 infections is desperately needed. This proposed EHR-based data mining offers a rapid
24 and empirically grounded approach to collecting such evidence and to informing the design of
25 prospective clinical trials that can focus on inflammatory pathways, biophysiological evidence
26 of the coinfection, and socio-behavioral determinants.
27
28
29
30
31
32
33
34

35 **Methods and Analysis**

36 ***Data Description***

37 We will use EHR from National COVID Cohort Collaborative (N3C). As of 8/2022, N3C has
38 aggregated 15.2 million patients (5.8 million COVID-19 patients) from 50 states¹¹. The EHR are
39 normalized by the Observational Health Data Sciences and Informatics (OHDSI)'s Observational
40 Medical Outcomes Partnership (OMOP) Common Data Model (CDM)¹². EHR variables
41 (individual level) span demographics, encounters, medical/social history, diagnoses,
42 procedures, medication prescriptions, medication exposure (e.g., vaccines), laboratory
43 tests/results, etc. Clinical notes have already been annotated in the CDM¹³. N3C has
44 epidemiological and community data (population level) including temporal COVID-19 burden,
45 vaccination, viral variance, health systems data, and geospatial data. N3C has pre-pandemic
46 data since 2018 and peri-pandemic data to the present. All data were de-identified and
47 updated every 2 weeks. We have the highest level of data access, which allows using patients'
48 residential ZIP codes and dates of clinical events. As of August of 2021, N3C has at least 13
49 thousand adults have been diagnosed and/or have had laboratory-confirmed HIV. A pilot study
50 shows that patients with coinfection have a high risk of hospitalization and mortality as
51 compared to COVID-19 patients without HIV infection¹⁴.
52
53
54
55
56
57
58
59
60

EHR Data Modeling

We will remodel the EHR data extracted from N3C. Because Individuals' longitudinal health records are stored at distributed locations in EHR. Many clinical events do not have explicit and/or complete temporal information. Raw EHR data as such are of little value for understanding questions such as why individuals with certain coinfection present different clinical outcomes and disease prognosis.

We will locate relevant phenotypic traits from N3C by curating OMOP CDM concept sets, a procedure called electronic phenotyping. Using these concept sets, we will then retrieve and integrate individuals' phenotypic traits from N3C. At last, we will retrieve temporal information for clinical events and establish a graphical model¹⁵ to represent clinical events and their temporal information.

EHR phenotyping. Phenotyping is the process of identifying cohorts and variables from raw EHR¹⁶. Because we use OMOP CDM-normalized EHR, phenotyping is the process of finding "OMOP concepts" that correspond to specific cohorts (e.g., all patients with ART-related visits) or variables (e.g., myocardial infarction)¹⁷. An OMOP concept is a unique identifier that is mapped to diverse medical codes that have the same semantic meaning but may be from different medical nomenclatures and EHR systems.

We will use standard phenotyping procedures. The logical procedures are as follows: 1) To identify existing OMOP concept sets (available in OHDSI's Atlas system and N3C) that can be used with minor revision. 2) If no appropriate concept sets available, we will follow the generic phenotyping procedures^{16,18} and curate OMOP CDM concepts from the Athena vocabulary repository, which allows retrieval of OMOP CDM concepts¹⁹. 3) To validate the revised or newly developed concept sets by using EHR chart review that is performed independently by two domain experts¹⁶.

Data retrieval and integration. To link external epidemiological and community-level data with individual-level EHR data, we will use individual patients' residential ZIP/county as the reference. We will use data imputation algorithms to impute missing values because population-towards-individual data integration automatically creates missing values. For existing missing values in EHR, we will infer missing values based on semantic relationships of OMOP CDM concepts. For geographical locations, we will use cities, counties, and states to infer locations at appropriate levels. For the rest of the missing values, we will use multiple imputation methods. Specifically, we will use selection models or pattern-mixture models for systematic missingness. We will also selectively use mean/median imputation, principal component analysis (PCA), singular value decomposition (SVD), k-nearest neighbor (kNN), least squares, expectation-maximization, and random forest.^{20,21} Data retrieval and integration will be implemented using SQL, R, Python, and PySpark, whichever appropriate.

Graphical EHR model. We will develop a customized Graphical model to represent the temporal relations among patients' demographic, clinical, and socio-behavioral data. A Graphical model encodes clinical events as nodes and their semantic relations (including temporal relations) as edges. **Figure 1** shows the proposed design of the model. General modeling procedures include: First, we will retrieve time stamps of clinical events. Some clinical events have explicit time stamps as captured by structured EHR. Many others do not (e.g., patient-reported symptoms as documented in clinical notes, trimester and/or gestational age). For those that do not have explicit time stamps, we will infer the time of the event based on neighboring EHR data. For

1
2
3 example, symptom onsets could be found in clinical notes (e.g., admission, discharge); trimester
4 and gestational age could be estimated by gestation-related diagnoses (e.g., Z3A: weeks of
5 gestation), procedures (e.g., ultrasound procedures), and clinical notes (e.g., last menstrual
6 period)^{22,23}. Second, we will represent temporal information of clinical events by modeling the
7 occurrences of an event and the instantaneous impact of the event. The occurrences of an
8 event are resulted from the first step. The instantaneous impact of an event is formulated using
9 exponential kernel functions and association rules based on clinical observation. Intuitively, two
10 clinical events with a long interval in between would have less impact on one another, but this
11 can be overwritten by events that hold special clinical meanings. Third, we will create recurrent
12 states for every clinical event by embedding the event, its semantic relations (i.e., edge in a
13 Graph) including the instantaneous impact of an event, and time. These recurrent states will
14 form a recurrent layer to be used for training machine learning models, which will be discussed
15 later.
16
17
18
19

20 **Machine Learning Modeling**

21 We will use supervised machine learning to examine patterns and sequences of phenotypic
22 traits for their predictive ability in clinical outcomes and disease prognosis. Existing studies
23 conclude differently on clinical outcomes among individuals with coinfection as well as the
24 factors correlated with these clinical outcomes. We provide two hypotheses. Hypothesis 1:
25 Individual patients respond differently to the coinfection. Hypothesis 2: Patients' clinical
26 outcomes and disease prognosis are attributed by the temporal dynamics of clinical events. If
27 these hypotheses are successfully tested, we will be able to delineate the impact of coinfection
28 on individuals' clinical outcomes. Therefore, we will customize recurrent neural network (RNN)
29 models to be used for predicting clinical outcomes and disease prognosis in real-time by
30 learning about patients' retrospective EHR at the individual level (personalized) as time
31 progresses. We adopt RNN for its unique advantage in capturing temporal dependencies of
32 data²⁴. Trained RNN models will be tested for their performance where the best-performed
33 model will be identified for identifying patterns/sequences of phenotypic traits predictive of
34 clinical outcomes and prognosis.
35
36
37

38 Cohort. Based on the estimated >13 thousand patients with the coinfection in our dataset, we
39 will blend in controls (COVID-19 patients without HIV) for each output variable using a match
40 ratio of 1:2, stratified by sex, race/ethnicity, and age. For a possible occasions of small sample,
41 for example, death cases (n<1000 with coinfection), the alternative strategy is to create
42 synthetic cases to impute and balance the sample.
43

44 Machine learning input. A complete and longitudinal health history together with linked
45 external epidemiological and community-level data will be included as the input of machine
46 learning models by which the models can learn from the input to predict individuals' in-time
47 clinical outcomes and disease prognosis. We will include but are not limited to the following
48 phenotypic traits: demographics, SDOH, diagnoses, underlying conditions, vitals, laboratory
49 tests, procedures, medication prescriptions/dispensing, medication exposure (e.g., vaccine),
50 and annotated clinical notes. Because there is no gold standard measure for ART adherence, we
51 will use "multiple measures" to estimate levels of ART adherence²⁵. Multiple measures include
52 medication events (inpatient dispensing), HIV-1 RNA copies (laboratory results), and medication
53 adherence data from clinical notes. For those with complete medication adherence data we use
54
55
56
57
58
59
60

1
2
3 the proportion of days covered (PDC) to categorize ART adherence levels (e.g., < 50, 50-80, 80-
4 85, 85-90, ≥90%)²⁶. We will categorize antiviral medications (ARV) into Integrase inhibitor
5 (INSTI)-based, non-nucleoside reverse transcriptase inhibitor (NNRTI)-based, protease inhibitor
6 (PI)-based, and other regimens. The approach to measuring ART adherence using EHR has
7 limitations, but the limitations can be mitigated by the well-presented and large-scale national
8 data. We collect both CD4⁺ counts as an indicator of existing damage and plasma HIV-1 RNA
9 copies as an indicator of projected disease progression.

10
11
12 ***Machine learning output.*** Clinical outcome measures as machine learning output include
13 inpatient admissions, length of stay (LOS), ICU admission, ICU LOS, comorbidities, and primary
14 discharge diagnosis²⁷. In addition to the measures within the acute phase of COVID-19, we will
15 also use symptoms, diagnoses, comorbidities, and PASC-associated-readmissions as outcome
16 and prognosis measures for individuals in the post-acute phase.

17
18 ***Model design.*** We will use RNN as the machine learning architecture to learn from patients'
19 longitudinal EHR and make the prediction of current and future clinical outcomes and disease
20 prognosis. We adopt RNN models because this neural network architecture is specialized for
21 capturing temporal dependency among event sequences. The RNN models will be trained to
22 learn from model input and to make the prediction of model output. With respect to the
23 embedding, we will include standard long short-term memory (LSTM) as well as phased LSTM
24 and other variants wherever appropriate²⁴. We will use the bag-of-pattern matrix as the
25 baseline embedding method to be compared against LSTM, in which this baseline method does
26 not fully consider temporal dependency. To test against RNN, we will use Support Vector
27 Machine (SVM) as the benchmark algorithm, in which SVM is a well-performed kernel-based
28 algorithm²⁸ but does not take full advantage of temporal dependency (hypothesis 2) and
29 personalized health records (hypothesis 1). Because the nature of our machine learning output
30 is binary variables, the proposed machine learning tasks are essentially binary classification
31 tasks. We will use Python for machine learning modeling.

32
33
34
35 ***Model evaluation (internal validity).*** To test the effectiveness of the prediction model, we will
36 use 10-fold cross validation. With respect to evaluation metrics, we will use the F score,
37 precision, recall, and the Area Under the Receiver Operating Characteristic (AUC) to assess the
38 models' predictive performance. We expect the F score, assuming balanced data, to reach a
39 minimum of 0.8. If trained models fail to meet this expectation, alternative strategies include
40 manually adding features handpicked by researchers after error analysis of models.

41 42 43 ***Clinical Decision Support System***

44
45 Based on the automatic clinical outcomes and disease prognosis prediction model, we will
46 design and implement a clinical decision support (CDS) prototype in collaboration with Prisma
47 Health clinics. The proposed CDS prototype is anticipated to assist providers in screening and
48 identifying patients who are at high risk of worse COVID-19 clinical outcomes (see **Table 1**), and
49 worse disease prognosis, including individuals with PASC. Specifically, the CDS will identify
50 individuals with a high risk of disease progression from their current clinical state (e.g., not
51 hospitalized, hospitalized, post-acute phase) by learning from the trained machine-learning
52 models. The effectiveness of the CDS demonstrates the external validity of the internally
53 validated predictive model and will be assessed by 1) appropriate identification for at-risk
54 individuals, 2) appropriate clinical actions, and 3) CDS system usability.

Table 1. Patient state according to WHO clinical progression scale.

Patient state	Clinical characteristics	Severity score
Uninfected	Uninfected, no viral RNA detected	0
Ambulatory mild disease	Asymptomatic, viral RNA detected	1
	Symptomatic, independent	2
	Symptomatic, assistance needed	3
Hospitalized, moderate disease	Hospitalized, no oxygen therapy	4
	Hospitalized, oxygen by mask or nasal prongs	5
	Hospitalized, oxygen by NIV or high flow	6
Hospitalized, severe disease	Intubation and mechanical ventilation, $pO_2/FiO_2 \geq 150$ or $SpO_2/FiO_2 \geq 200$	7
	Mechanical ventilation, $pO_2/FiO_2 < 150$ ($SpO_2/FiO_2 < 200$) or vasopressors	8
	Mechanical ventilation, $pO_2/FiO_2 < 150$ and vasopressors, dialysis, or ECMO	9
Dead	Dead	10

CDS workflow. The proposed CDS is a hybrid of knowledge-based and non-knowledge-based system²⁹. It has 1) a machine-learning-based prediction module (non-knowledge-based) for identifying high-risk patients and 2) a provider-curated medical logic module (knowledge-based) for generating clinical actions for identified high-risk patients. The CDS testing takes place in a retrospective way (i.e., using retrospective EHR).

Cohort definition and data collection. Using retrospective EHR data (2-year baseline~2023) from Prisma Health's Epic system, we will first group the existing PLWH who have COVID-19 (sampling $n > 500$) based on their state at the point of CDS screening. The patient states include 1) ambulatory patients with COVID-19, 2) hospitalized patients for COVID-19 with moderate disease, 3) hospitalized patients for COVID-19 with severe disease, and 4) post-acute phase of COVID-19 (i.e., from beyond 4 weeks after symptom onset)³⁰. See **Table 1** for definitions of states 1-3 based on WHO's clinical progression scale for COVID-19.

Prediction module. For patients in each state, we will use the trained machine-learning model to learn from previous medical records and predict worsening clinical outcomes as time progresses (i.e., acute, and post-acute phases every 3 months). The prediction will include primary COVID-19 clinical outcomes (**Table 2**) developed by the WHO Working Group on the Clinical Characterisation and Management of COVID-19³¹.

Table 2. Key clinical outcome measures.

Organ dysfunction
• Murray score
• Sequential organ failure assessment score, multiple organ dysfunction score
• Acute coronary syndrome; arrhythmias
• Delirium

Comorbidities <ul style="list-style-type: none"> • Pulmonary, cardiovascular, renal, neurological, etc.
Secondary infection <ul style="list-style-type: none"> • Bacterial, viral
Biochemical parameters <ul style="list-style-type: none"> • C-reactive protein, D-dimers, IL-6, and ferritin serum concentrations, and leucocyte counts
Radiological findings <ul style="list-style-type: none"> • Chest CT scan, chest x-ray
Duration of intervention <ul style="list-style-type: none"> • Inpatient admission, length of stay (LOS) • ICU admission, ICU LOS • Ventilation • Organ support or hospital-free days
Pregnancy outcomes <ul style="list-style-type: none"> • Preterm delivery, miscarriage • Fetal status • Severe maternal morbidity (SMM) measures
Mortality <ul style="list-style-type: none"> • All-cause mortality at hospital discharge
Quality of life <ul style="list-style-type: none"> • Longer term survival and primary diagnoses for readmission (post-acute phase)

Medical logic module. Patients identified by the CDS to have an increased risk of worse clinical outcomes will be reviewed and discussed by two providers who are specialized in HIV and COVID-19. First, the providers will generate gold-standard judgment on whether a patient is correctly identified by the prediction module, which later will be used for assessing the effectiveness of CDS. Second, the providers will generate appropriate clinical actions upon chart review. These clinical actions will be made up to date with the “NIH Guidance for COVID-19 and People with HIV”, including treatment options based on cohorts and risk factors, medication reconciliation considering ART regimens, consultation with specialists for multi-organ system complications and PASC, referrals, and outreach³². Providers’ decision-making processes will be programmed using Arden Syntax (v3) or Clinical Quality Language³³ in the knowledge base, which is determined by specific Epic EHR data model.

Effectiveness of CDS (external validity). There are two evaluation metrics: 1) Appropriate identification for individuals at high risk for adverse clinical outcomes (**Table 2**) by comparing model-identified cases against the gold standard generated from chart review. We will use F measure (>0.8), AUC, precision, and recall for assessment; 2) Appropriate clinical actions using a quasi-experimental design. We will compare outcomes of patients who naturally used the medical logic module-suggested care against those who did not (n=100 each). The outcomes include but are not limited to readmissions (e.g., same day, 7-, 14-, 30- days), healthcare utilization (e.g., LOS, ER/observation visits, ICU admission). We will use mixed-effect generalized regression models to estimate model effectiveness wherever appropriate.

1
2
3 *Usability testing.* We will assess CDS usability by adopting the “think aloud” protocol³⁴. The two
4 providers from Prisma Health will participate in the test. Each one will be presented with
5 randomly selected EHR (n=5 at-risk cases + n=5 control cases) along with the CDS output. In
6 each case, participants will be instructed to verbalize their reasoning procedures (e.g.,
7 phenotypic traits from EHR that can be used in the reasoning, logic flow) towards identifying at-
8 risk patients and corresponding clinical decisions. Sessions are audio-recorded and will then be
9 coded (e.g., by content, understandability, navigation, workflow, visibility, and usability)
10 independently by two researchers for downstream analyses.
11
12
13

14 **Patient and Public Involvement**

15 No patient involved.
16
17
18

19 **Ethics and Dissemination**

20 The study was approved by the institutional review boards at the University of South Carolina
21 (Pro00121828) as non-human subject study.
22

23 This study will result in a comprehensive knowledge base that documents clinical outcomes and
24 disease prognosis for individuals with the coinfection, their risk factors (e.g., underlying
25 conditions, ART adherence, comorbidities, socio-behavioral), and their responses to
26 therapeutics. This study will also result in a prototype CDS that can identify patients at high risk
27 of worsening clinical outcomes and prognosis in real-time. These results are generalizable and
28 will form a foundation for developing comprehensive real-world CDS systems for
29 implementation in state-wide and national HIV and COVID-19 clinics.
30

31 Study findings will be presented at academic conferences and published in peer-reviewed
32 journals. This study will disseminate urgently needed clinical evidence for guiding clinical
33 practice for individuals with the coinfection at Prisma Health.
34
35
36
37

38 **Authors' Contributions**

39 CL conceived the study design and drafted the manuscript. CL completed preliminary data
40 collection. SW, BO, EP, MY, and XL contributed critical edits to the manuscript. All authors
41 reviewed and approved the manuscript.
42
43
44

45 **Funding Statement**

46 Research reported in this publication was supported by the National Institute Of Allergy And
47 Infectious Diseases of the National Institutes of Health under Award Number R21AI170171. The
48 content is solely the responsibility of the authors and does not necessarily represent the official
49 views of the National Institutes of Health.
50
51
52

53 **Competing Interest Statement**

54 None declared.
55
56
57
58
59
60

References

1. World Health Organization. *Clinical Features and Prognostic Factors of COVID-19 in People Living with HIV Hospitalized with Suspected or Confirmed SARS-CoV-2 Infection.*; 2021.
2. Dandachi D, Geiger G, Montgomery MW, et al. Characteristics, Comorbidities, and outcomes in a multicenter registry of patients with human immunodeficiency virus and coronavirus disease 2019. *Clin Infect Dis*. Published online 2020.
3. Braunstein SL, Lazar R, Wahnich A, Daskalakis DC, Blackstock OJ. COVID-19 infection among people with HIV in New York City: A population-level analysis of linked surveillance data. *Clinical Infectious Diseases: An Official Publication of the Infectious Diseases Society of America*. Published online 2020.
4. Eisinger RW, Lerner AM, Fauci AS. Human Immunodeficiency Virus/AIDS in the Era of Coronavirus Disease 2019: A Juxtaposition of 2 Pandemics. *The Journal of Infectious Diseases*. Published online 2021.
5. Cooper TJ, Woodward BL, Alom S, Harky A. Coronavirus disease 2019 (COVID-19) outcomes in HIV/AIDS patients: a systematic review. *HIV Med*. 2020;21(9):567-577.
6. Calza L, Bon I, Tadolini M, et al. COVID-19 in patients with HIV-1 infection: a single-centre experience in northern Italy. *Infection*. 2021;49(2):333-337.
7. Costenaro P, Minotti C, Barbieri E, Giaquinto C, Donà D. SARS-CoV-2 infection in people living with HIV: a systematic review. *Reviews in Medical Virology*. 2021;31(1):1-12.
8. Park LS, Rentsch CT, Sigel K, et al. COVID-19 in the largest US HIV cohort. *AIDS*. 2020;2020:23rd.
9. Haendel MA, Chute CG, Gersing K. The National COVID Cohort Collaborative (N3C): rationale, design, infrastructure, and deployment. *Journal of the American Medical Informatics Association*. Published online 2020.
10. Deer RR, Rock MA, Vasilevsky N, et al. Characterizing long COVID: deep phenotype of a complex condition. *EBioMedicine*. 2021;74:103722.
11. Datavent. COVID-19 Research Database. Accessed February 20, 2021. <https://covid19researchdatabase.org/>
12. OHDSI community. Observational Health Data Sciences and Informatics Common Data Model.
13. N3C. COVID-19 Clinical Data Warehouse Data Dictionary.
14. Yang X, Zhang J, Guo S, Olatosi B, Weissman S, Li X. The role of HIV infection in the clinical spectrum of COVID-19: a population-based cohort analysis based on US National COVID Cohort Collaborative (N3C) Enclave data. *Available at SSRN 3860395*. Published online 2021.
15. Liu C, Wang F, Hu J, Xiong H. Temporal phenotyping from longitudinal electronic health records: A graph based framework. In: *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining.* ; 2015:705-714.
16. Banda JM, Seneviratne M, Hernandez-Boussard T, Shah NH. Advances in electronic phenotyping: from rule-based definitions to machine learning models. *Annu Rev Biomed Data Sci*. 2018;1:53-68.
17. Richesson RL, Hammond WE, Nahm M, et al. Electronic health records based phenotyping in next-generation clinical trials: a perspective from the NIH Health Care

- Systems Collaboratory. *Journal of the American Medical Informatics Association*. 2013;20(e2):e226-e231.
18. Weng C, Shah NH, Hripcsak G. Deep phenotyping: Embracing complexity and temporality—Towards scalability, portability, and interoperability. *Journal of Biomedical Informatics*. 2020;105:103433.
 19. OHDSI Athena standard vocabularies. Accessed September 1, 2021. <https://www.ohdsi.org/analytic-tools/athena-standardized-vocabularies/>
 20. Luo Y. Evaluating the state of the art in missing data imputation for clinical data. *Briefings in Bioinformatics*. 2022;23(1):bbab489.
 21. Li J, Yan XS, Chaudhary D, et al. Imputation of missing values for electronic health record laboratory data. *NPJ Digit Med*. 2021;4(1):1-14.
 22. Lyu T, Liang C, Liu J, et al. Temporal Events Detector for Pregnancy Care (TED-PC): A Rule-based Algorithm to Infer Gestational Age and Delivery Date from Electronic Health Records of Pregnant Women with and without COVID-19. *arXiv preprint arXiv:220502933*. Published online 2022.
 23. Liu J, Hung P, Liang C, et al. Multilevel determinants of racial/ethnic disparities in severe maternal morbidity and mortality in the context of the COVID-19 pandemic in the USA: protocol for a concurrent triangulation, mixed-methods study. *BMJ Open*. 2022;12(6):e062294.
 24. Goodfellow I, Bengio Y, Courville A, Bengio Y. *Deep Learning*. Vol 1. MIT press Cambridge; 2016.
 25. Castillo-Mancilla JR, Haberer JE. Adherence measurements in HIV: new advancements in pharmacologic methods and real-time monitoring. *Current Hiv/aids Reports*. 2018;15(1):49-59.
 26. Byrd KK, Hou JG, Hazen R, et al. Antiretroviral adherence level necessary for HIV viral suppression using real-world data. *J Acquir Immune Defic Syndr*. 2019;82(3):245.
 27. Lavery AM, Preston LE, Ko JY, et al. Characteristics of Hospitalized COVID-19 Patients Discharged and Experiencing Same-Hospital Readmission—United States, March--August 2020. *Morbidity and Mortality Weekly Report*. 2020;69(45):1695.
 28. Murphy KP. *Machine Learning: A Probabilistic Perspective*. MIT press; 2012.
 29. Shiffman RN, Wright A. Evidence-based clinical decision support. *Yearb Med Inform*. 2013;22(01):120-127.
 30. Nalbandian A, Sehgal K, Gupta A, et al. Post-acute COVID-19 syndrome. *Nat Med*. Published online 2021:1-15.
 31. Marshall JC, Murthy S, Diaz J, et al. A minimal common outcome measure set for COVID-19 clinical research. *The Lancet Infectious Diseases*. 2020;20(8):e192–e197.
 32. Guidelines Working Groups of the NIH Office of AIDS Research Advisory Council. Guidance for COVID-19 and People with HIV.
 33. Hripcsak G, Clayton P, Pryor T, Haug P, Wigertz O, der Lei J. The Arden syntax for medical logic modules. In: *Proceedings. Symposium on Computer Applications in Medical Care*. ; 1990:200-204.
 34. Li AC, Kannry JL, Kushniruk A, et al. Integrating usability testing and think-aloud protocol analysis with “near-live” clinical simulations in evaluating clinical decision support. *Int J Med Inform*. 2012;81(11):761-772.

For peer review only

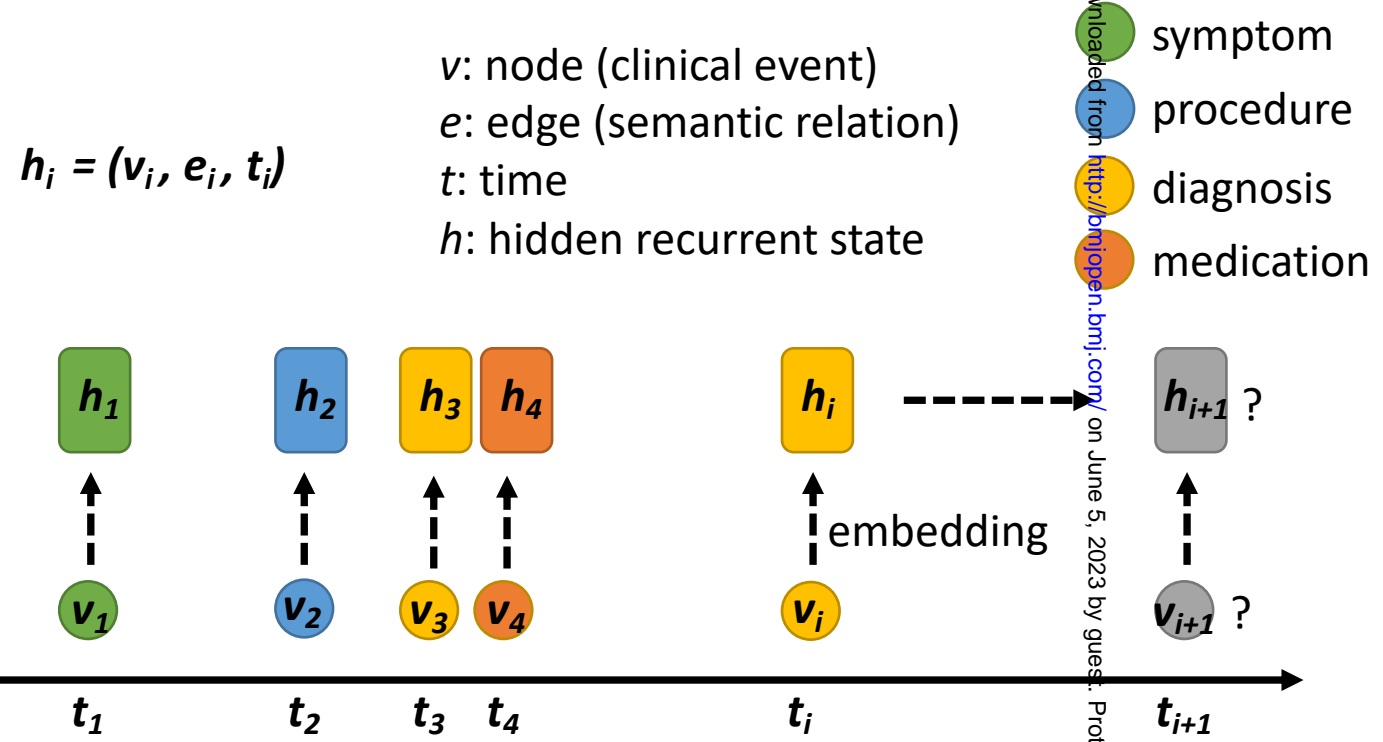
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

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

Figure 1. EHR model design.

For peer review only

067204 on 13 September 2022. Downloaded from <http://bmjopen.bmj.com/> on June 5, 2023 by guest. Protected by copyright.





CONSORT 2010 checklist of information to include when reporting a randomised trial*

Section/Topic	Item No	Checklist item	Reported on page No
Title and abstract			
	1a	Identification as a randomised trial in the title	NA
	1b	Structured summary of trial design, methods, results, and conclusions (for specific guidance see CONSORT for abstracts)	2
Introduction			
Background and objectives	2a	Scientific background and explanation of rationale	4-5
	2b	Specific objectives or hypotheses	4-5
Methods			
Trial design	3a	Description of trial design (such as parallel, factorial) including allocation ratio	NA
	3b	Important changes to methods after trial commencement (such as eligibility criteria), with reasons	NA
Participants	4a	Eligibility criteria for participants	5, 8
	4b	Settings and locations where the data were collected	5
Interventions	5	The interventions for each group with sufficient details to allow replication, including how and when they were actually administered	NA
Outcomes	6a	Completely defined pre-specified primary and secondary outcome measures, including how and when they were assessed	11
	6b	Any changes to trial outcomes after the trial commenced, with reasons	NA
Sample size	7a	How sample size was determined	8
	7b	When applicable, explanation of any interim analyses and stopping guidelines	NA
Randomisation:			
Sequence generation	8a	Method used to generate the random allocation sequence	NA
	8b	Type of randomisation; details of any restriction (such as blocking and block size)	NA
Allocation concealment mechanism	9	Mechanism used to implement the random allocation sequence (such as sequentially numbered containers), describing any steps taken to conceal the sequence until interventions were assigned	NA
Implementation	10	Who generated the random allocation sequence, who enrolled participants, and who assigned participants to interventions	NA
Blinding	11a	If done, who was blinded after assignment to interventions (for example, participants, care providers, those	NA

		assessing outcomes) and how	
	11b	If relevant, description of the similarity of interventions	NA
Statistical methods	12a	Statistical methods used to compare groups for primary and secondary outcomes	NA
	12b	Methods for additional analyses, such as subgroup analyses and adjusted analyses	NA
Results			
Participant flow (a diagram is strongly recommended)	13a	For each group, the numbers of participants who were randomly assigned, received intended treatment, and were analysed for the primary outcome	NA
	13b	For each group, losses and exclusions after randomisation, together with reasons	NA
Recruitment	14a	Dates defining the periods of recruitment and follow-up	NA
	14b	Why the trial ended or was stopped	NA
Baseline data	15	A table showing baseline demographic and clinical characteristics for each group	NA
Numbers analysed	16	For each group, number of participants (denominator) included in each analysis and whether the analysis was by original assigned groups	8
Outcomes and estimation	17a	For each primary and secondary outcome, results for each group, and the estimated effect size and its precision (such as 95% confidence interval)	NA
	17b	For binary outcomes, presentation of both absolute and relative effect sizes is recommended	NA
Ancillary analyses	18	Results of any other analyses performed, including subgroup analyses and adjusted analyses, distinguishing pre-specified from exploratory	11
Harms	19	All important harms or unintended effects in each group (for specific guidance see CONSORT for harms)	NA
Discussion			
Limitations	20	Trial limitations, addressing sources of potential bias, imprecision, and, if relevant, multiplicity of analyses	NA
Generalisability	21	Generalisability (external validity, applicability) of the trial findings	9-11
Interpretation	22	Interpretation consistent with results, balancing benefits and harms, and considering other relevant evidence	9-11
Other information			
Registration	23	Registration number and name of trial registry	NA
Protocol	24	Where the full trial protocol can be accessed, if available	NA
Funding	25	Sources of funding and other support (such as supply of drugs), role of funders	12

*We strongly recommend reading this statement in conjunction with the CONSORT 2010 Explanation and Elaboration for important clarifications on all the items. If relevant, we also recommend reading CONSORT extensions for cluster randomised trials, non-inferiority and equivalence trials, non-pharmacological treatments, herbal interventions, and pragmatic trials. Additional extensions are forthcoming: for those and for up to date references relevant to this checklist, see www.consort-statement.org.