

ISPE-Endorsed Guidance in Using Electronic Health Records for Comparative Effectiveness Research in COVID-19: Opportunities and Trade-Offs

Grammati Sarri^{1,*}, Dimitri Bennett^{2,3}, Thomas Debray^{4,5}, Anouk Deruaz-Luyet⁶, Montse Soriano Gabarró⁷, Joan A. Largent⁸, Xiaojuan Li⁹, Wei Liu¹⁰, Jennifer L. Lund¹¹, Daniela C. Moga¹², Mugdha Gokhale^{11,13}, Christopher T. Rentsch^{14,15}, Xuerong Wen¹⁶, Chen Yanover¹⁷, Yizhou Ye¹⁸, Huifeng Yun¹⁹, Andrew R. Zullo^{20,21,22,23} and Kueiyu Joshua Lin²⁴

As the scientific research community along with healthcare professionals and decision makers around the world fight tirelessly against the coronavirus disease 2019 (COVID-19) pandemic, the need for comparative effectiveness research (CER) on preventive and therapeutic interventions for COVID-19 is immense. Randomized controlled trials markedly under-represent the frail and complex patients seen in routine care, and they do not typically have data on long-term treatment effects. The increasing availability of electronic health records (EHRs) for clinical research offers the opportunity to generate timely real-world evidence reflective of routine care for optimal management of COVID-19. However, there are many potential threats to the validity of CER based on EHR data that are not originally generated for research purposes. To ensure unbiased and robust results, we need high-quality healthcare databases, rigorous study designs, and proper implementation of appropriate statistical methods. We aimed to describe opportunities and challenges in EHR-based CER for COVID-19-related questions and to introduce best practices in pharmacoepidemiology to minimize potential biases. We structured our discussion into the following topics: (1) study population identification based on exposure status; (2) ascertainment of outcomes; (3) common biases and potential solutions; and (iv) data operational challenges specific to COVID-19 CER using EHRs. We provide structured guidance for the proper conduct and appraisal of drug and vaccine effectiveness and safety research using EHR data for the pandemic. This paper is endorsed by the International Society for Pharmacoepidemiology (ISPE).

Coronavirus disease 2019 (COVID-19), caused by the severe acute respiratory syndrome coronavirus-2 (SARS-CoV-2), emerged in 2019 as a major and urgent public health emergency worldwide.¹ With the number of known cases and deaths rising exponentially (as of January 2022, there were over 356 million confirmed cases and over 5.6 million deaths),² public health

control measures have focused on improving preventive strategies, including the introduction of nonpharmaceutical interventions, improving testing facilities, and restrictive social measures. Declared as a pandemic on March 11, 2020, COVID-19 unavoidably continues to place a huge strain on our activities of daily living while posing significant health, social, economic,

¹Visiting Lead Scientist, Cytel Inc., London, UK; ²Takeda Global Evidence and Outcomes, Takeda Pharmaceuticals USA, Inc., Cambridge, Massachusetts, USA; ³Center for Clinical Epidemiology and Biostatistics, Perelman School of Medicine at the University of Pennsylvania, Philadelphia, Pennsylvania, USA; ⁴Julius Center for Health Sciences and Primary Care, University Medical Centre Utrecht, Utrecht, The Netherlands; ⁵Smart Data Analysis and Statistics, Utrecht, The Netherlands; ⁶Global Epidemiology and Real-World Evidence CoE, Corporate Medical Affairs, Boehringer Ingelheim International GmbH, Ingelheim-am-Rhein, Germany; ⁷Bayer Partnerships and Integrated Evidence Generation Office, Integrated Evidence Generation & Business Innovation, Medical Affairs & Pharmacovigilance, Bayer AG, Berlin, Germany; ⁸Real-World Solutions, IQVIA, Mission Viejo, California, USA; ⁹Department of Population Medicine, Harvard Medical School and Harvard Pilgrim Health Care Institute, Boston, Massachusetts, USA; ¹⁰Division of Epidemiology, Office of Surveillance and Epidemiology, Center for Drug Evaluation and Research, Food and Drug Administration, Silver Spring, Maryland, USA; ¹¹Department of Epidemiology, Gillings School of Global Public Health, University of North Carolina at Chapel Hill, Chapel Hill, North Carolina, USA; ¹²Department of Pharmacy Practice and Science, College of Pharmacy, University of Kentucky, Lexington, Kentucky, USA; ¹³Department of Epidemiology, Merck, West Point, Pennsylvania, USA; ¹⁴Faculty of Epidemiology and Population Health, Department of Non-communicable Disease Epidemiology, London School of Hygiene & Tropical Medicine, London, UK; ¹⁵Department of Internal Medicine, Yale School of Medicine, New Haven, Connecticut, USA; ¹⁶Health Outcomes, Pharmacy Practice, College of Pharmacy, University of Rhode Island, Kingston, Rhode Island, USA; ¹⁷KI Research Institute, Kfar Malal, Israel; ¹⁸Global Epidemiology, Pharmacovigilance and Patient Safety, AbbVie Inc., North Chicago, Illinois, USA; ¹⁹Department of Epidemiology, University of Alabama at Birmingham, Birmingham, Alabama, USA; ²⁰Department of Health Services, Policy, and Practice, Brown University School of Public Health, Providence, Rhode Island, USA; ²¹Department of Epidemiology, Brown University School of Public Health, Providence, Rhode Island, USA; ²²Center of Innovation in Long-Term Services and Supports, Providence Veterans Affairs Medical Center, Providence, Rhode Island, USA; ²³Department of Pharmacy, Lifespan-Rhode Island Hospital, Providence, Rhode Island, USA; ²⁴Brigham and Women's Hospital and Harvard Medical School, Boston, Massachusetts, USA. *Correspondence: Grammati Sarri (grammati.sarri@gmail.com)

Received November 20, 2021; accepted February 2, 2022. doi:10.1002/cpt.2560

and environmental challenges with major implications for the entire global community.³

The scientific research community has tirelessly worked on the fight against the virus, the disease, and its complications. COVID-19 vaccines have been introduced^{4–6} and therapies are being developed and studied.⁷ Whereas large-scale randomized controlled trials (RCTs) remain pivotal to determine an intervention's efficacy and safety for regulatory approval purposes,⁸ many medications and vaccines have received emergency use authorizations (EUAs) or fast-tracked approval, which resulted in their extensive use in wider populations (beyond the targeted populations recruited in RCTs). In addition, some primary concerns were raised about the use of RCTs to generate generalizable evidence in COVID-19. The selected participants in the RCTs may be unlikely to represent the frail and complex patients seen in routine care, and the short duration of RCTs did not allow for the generation of findings on long-term safety and effectiveness outcomes. Given this, we need to also rely on nonrandomized studies to generate real-world evidence (RWE) on the effectiveness and safety of preventive and therapeutic interventions. RWE can be generated through a range of applications, like individualized prescribing, postmarketing surveillance, and can support policy or reimbursement decisions. Databases that reflect routine care delivery include electronic health records (EHRs), administrative health insurance claims, disease and product registries, and other non-research-specific data sources (for example, social media).⁹ The advantages of RWE providing timely, generalizable evidence from a large, diverse group of patients are well-recognized, particularly for a public health crisis like the COVID-19 pandemic.^{10,11} To capitalize on the RWE offerings to inform the management of COVID-19, we need access to high-quality healthcare data, rigorous study designs, and proper implementation of appropriate statistical methods to ensure unbiased and robust results.¹²

EHRs have been widely used throughout the pandemic to generate evidence for risk stratification of patients, prognostic and risk factor identification, natural disease history investigation, and outcomes of interest that may be used in comparative effectiveness research (CER),¹² many of which are not routinely available in claim databases. In addition to its timely availability for research purposes, EHRs capture information on key factors for patient phenotyping and confounding adjustment, including inpatient medication use, disease, and patient characteristics (such as vital signs, laboratory test, and imaging results, smoking status, body mass index, and code status). EHRs also provide clinical notes and reports that are often important for validation studies.¹³ Furthermore, the growing availability and utilization of EHRs from different populations across healthcare systems with federated data networks and multi-database infrastructure involving several countries are also contributing to increasing opportunities for urgent and critical CER evidence on COVID-19 treatments.¹⁴ On the contrary, the availability of insurance claims and registry data is typically lagged, which poses a barrier to timely availability of such evidence.

However, controversies around the validity of findings from previous studies using real-world data (RWD) assessing the effectiveness of COVID-19 treatments have sparked skepticism around the use of this evidence to inform clinical decision making

in the management of the pandemic.¹⁵ In addition to concerns with the veracity and appropriateness of specific data sources, some of these controversies arose from the inherent limitations of RWE that are not unique in COVID-19 CER, including data quality, missing data, and confounding bias.¹⁶ Some challenges are EHR-specific, such as misclassification of key information due to EHR data-discontinuity,¹⁷ converting unstructured free-text data into structured data,¹⁸ and harmonization of data across EHR systems in a multicenter study.¹⁹ Given the growing requirement for optimizing the potential of RWD for assessing the real-time effectiveness and safety of COVID-19 treatments, there is a compelling need for setting up clear guidance to ensure the results of these observational studies are reliable and valid for decision making. The purpose of this paper is to discuss the opportunities, unique challenges, and potential solutions when using EHR data for CER to inform the delivery of care in response to public health crises, such as the COVID-19 pandemic. We hope this will equip the readers with a nonexhaustive list of tools to implement and interpret quality RWE using EHRs in a pandemic setting.

METHODOLOGICAL APPROACH

A targeted literature review conducted in March 2021 (for full details, please see **Supplementary Material**) was used to inform this guidance paper along with the discussions held among the participants of the International Society for Pharmacoepidemiology (ISPE) Comparative Effectiveness Research (CER) Special Interest Group (SIG) working group.

The purpose of these discussions was to identify opportunities, challenges, and good research practices around the use of EHRs in COVID-19 CER to inform the content of this paper, which was structured in two areas: (1) methodological issues including how to define “exposure” and “outcomes” in COVID-19 CER, how to minimize confounding and information bias, and other related methodological issues; and (2) data operational challenges specific to COVID-19 CER. A summary of these considerations is presented in **Figure 1**. Some of these concerns may be more relevant or applicable for one type of EHRs over another (inpatient, outpatient EHRs, post-acute care, long-term care settings, or linked EHRs). It is important to note that our discussion focused on CER issues arising after the selection of EHR data sources that have been tested for validity and reliability by investigators, rather than issues related to fitness-for-use of data from EHRs. The urgency of the COVID-19 pandemic may impose a need to explore new data sources; therefore, its fitness-for-use for research purposes should be thoroughly investigated. We also encourage our readers to consider our paper alongside previous relevant guidance related to the design of nonrandomized studies, data collection, source validation, results reproducibility, how to reliably synthesize results from RCTs and nonrandomized studies and on general topics regarding the use of this evidence in CER (**Table S1**).

METHODOLOGICAL ISSUES USING EHRs IN COVID-19 CER Defining the study population based on preventive or treatment exposure

In CER, the study population is typically defined by use or non-use of specific preventive and therapeutic interventions. For

Considerations in EHRs for COVID-19 CER

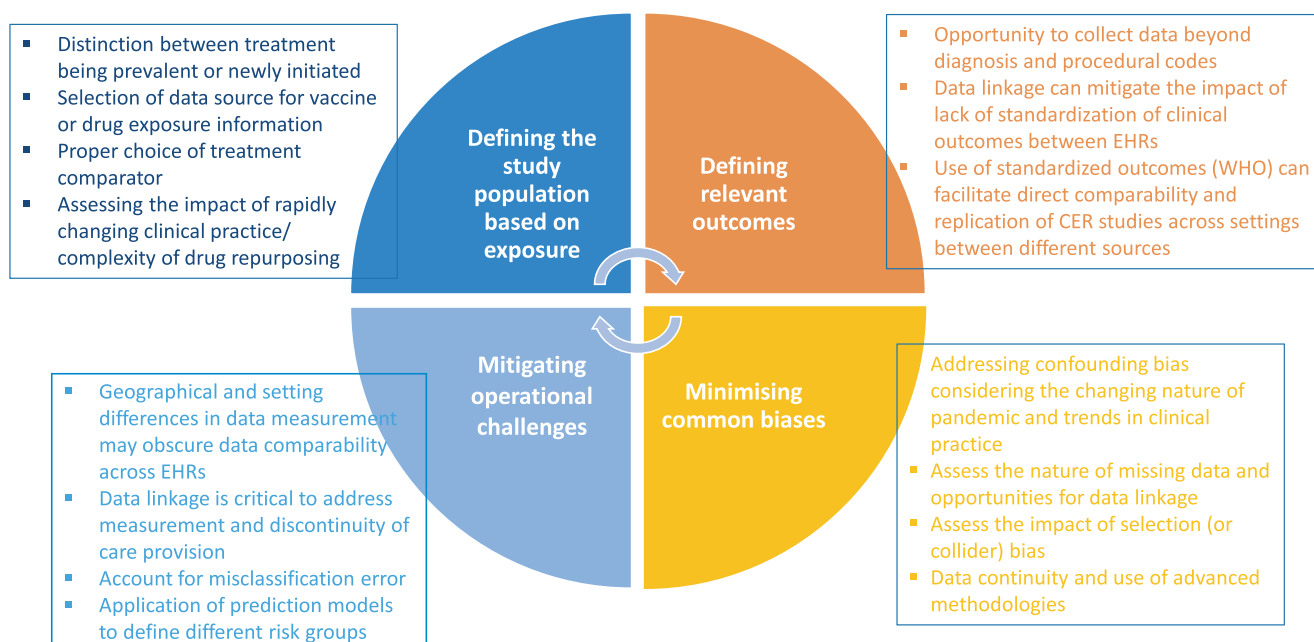


Figure 1 Considerations in EHRs for COVID-19 CER. CER, comparative effectiveness research; COVID-19, coronavirus disease 2019; EHRs, electronic health records; WHO, World Health Organization.

the purposes of this paper, exposure refers to pharmacological and nonpharmacological interventions used to prevent or treat COVID-19, including vaccines and therapeutics. EHR data typically allow researchers to define COVID-19 by positivity of laboratory results. We recommend using a case definition based on International Classification of Diseases (ICD) 10 diagnosis codes (e.g., U07.1) or positive results of a Nucleic Acid Amplification Test or reverse transcription polymerase chain reaction (PCR) because rapid antigen tests have variable performance when validated against the PCR test results, especially for asymptomatic patients or those with symptom onset more than 1 week ago.²⁰ Antibodies have not been routinely used for COVID-19 diagnosis and should be used in combination with other information (i.e., symptoms) to define COVID-19 in CER.²¹ The key considerations related to defining treatment exposure data using EHR data are as follows:

- **Selecting the appropriate source(s) for vaccine or drug exposure information:** EHR systems may have heterogeneous sources of medication information, each with various details and validity regarding preventive or therapeutic interventions. Prescribing (i.e., order) data are typically available in EHRs, and, in many settings, electronic medication administration record (eMAR) data as well. Some provider systems may have on-site pharmacies where dispensing data are also available. Medications details, such as dose, duration, and frequency, are usually available in EHRs, although researchers should be cautious on the validity of this information (e.g., frequency of the prescribing and dispensing data may be influenced by adherence, whereas eMARs recorded the actual timing of drug

administration). In general, the likelihood of exposure misclassification is the lowest based on eMARs among these data sources; dispensing data is less susceptible to information bias than prescribing data as the former is one step closer to actual ingestion of medication than the latter. The eMAR data have reliable inpatient, emergency department, or outpatient on-site drug use information. However, researchers should not only be concerned about ascertainment bias in relation to identification of cases but also assess data completeness for the specific drugs of interest and work with clinical experts within the system to identify these gaps as data completeness is contextual and is determined through an understanding of specific data needs.²²

In addition, it is important to note that vaccine information may be incompletely captured or missing in the EHRs as most COVID-19 vaccinations programs may occur in mass vaccination sites, pharmacies, and other settings where often no health insurance claims are submitted. Information on whether an individual has been vaccinated or not relies, at times, on the individual reporting such information to their healthcare professional. The propensity of misclassifying or missing information on the vaccination status is highly dependent on local vaccination policy settings. For example, in settings where vaccines are mandatory, there is an opportunity for data linkage to vaccination records or to the collection of patient-generated information, and such information could be captured in EHR-based vaccine studies. Moreover, in settings where significant under-recording of vaccination status may be an issue, a correction factor in outcomes analysis can be applied using a standard methodology to consider exposure misclassifications.²³

Alternatively, it is possible to adopt multiple imputation methods that recover missing vaccination status from observed information. Although imputation methods commonly assume that data are “missing at random,” several extensions have been proposed to account for more complex patterns of missingness.²⁴ For example, if unvaccinated people are less likely to provide information on vaccination status (leading to missingness is dependent on unobserved data situations (please see below for more information on missing bias issues)), Heckman selection models could be used to generate imputations under less stringent assumptions.²⁵ These models make use of instrumental variables that are not related to the exposure status (and therefore not informative in traditional imputation methods) but affect their registration (and thus availability). Common examples of instrumental variables are characteristics that describe (differences in) the data collection procedure, such as the method used to retrieve the vaccination status (e.g., interview, questionnaire, etc.).

EHR data also often include structured and unstructured data, and consideration of the source and accuracy of drug exposure information is important. Unstructured pharmacy data usually take the form of free-text fields in which providers record information about prescriptions, and therefore capture information with varying degrees of completeness. Although unstructured clinical notes and images may contain medications not available in prescribing or dispensing data, such as over-the-counter drug or supplement information, natural language processing (NLP) of the free-text notes is needed. NLP may also help extract medical indications, such as thrombotic events or bleeding risks that are contained in the unstructured EHR. However, developing a valid NLP module often requires manual chart review to establish the annotated dataset (the “gold-standard”), which is resource- and time-consuming.

- **Appropriate choice of treatment comparator:** Operationalization of an exposure-comparator definition must be mapped to the specified research question.²⁶ However, given the dynamic nature of COVID-19 (i.e., rapidly evolving knowledge of the disease’s natural history) and the lack of standard treatment guidelines, particularly at the early phase of the pandemic, it can be challenging to find appropriate therapeutic comparators. For example, each of the following comparisons may be faced with different methodological challenges: (a) comparing initiation of different treatment options (e.g., a drug or vaccine)²⁷; (b) comparing initiation of different doses of the same therapeutic agent; or (c) comparing any use vs. no use of a treatment or comparison of different sequential drug therapy strategies.²⁸ The use of an active-comparator, new-user study design is generally more desirable (comparisons (a) and (b)) as it aims to mitigate biases by first restricting the study to individuals with an indication for treatment and without contraindications, while also aligning individuals at the same point in time to start follow-up (i.e., treatment initiation) and ensuring the correct temporality between covariate and exposure assessment.²⁹ Researchers should pay close attention to the evolving guidelines on treatment options by disease severity (e.g., monoclonal antibodies are recommended to initiate early in the disease course,³⁰ and systemic steroids are indicated for patients

with moderate-to-severe disease³¹) while considering the use of active comparators of the same administration route.

Comparison between treatments given at different disease stages can lead to refractory confounding bias.³² Therefore, it is not surprising that researchers may experience difficulties in identifying a comparable alternative treatment in a newly emerging disease like COVID-19, which led to many nonuser comparisons. Nonuser comparisons in the nonrandomized settings could be subject to two types of bias: (1) immortal time bias: researchers often need an exposure assessment period to determine nonuse status (e.g., no use in the first 48 hours after hospital admission) as immortal time bias can occur if the immortal time (i.e., the person-time in which a patient “cannot die” by design) is differentially distributed among the two arms under comparison. This is often the case if the start of the follow-up began before the end of this exposure assessment period (e.g., the follow-up starts on the admission date) when comparing use vs. nonuse of a pharmacotherapy; and (2) confounding bias: when prescribing is highly informed by prognostic factors, nonusers are either much healthier individuals for whom no treatment is needed or are individuals with a grave prognosis for whom many aggressive treatments may be withheld. Such confounding may not be addressable if some prognostic factors are unmeasured in the study database. In addition, because COVID-19 care is rapidly changing as research findings emerge, careful consideration of the time trend of clinical practice is critical when choosing an appropriate comparator.³³

- **The distinction between prevalent vs. newly initiated users:** A new-user design is recommended in CER because the hazards of medical treatment may be different for a new user compared with a chronic user who had tolerated it before cohort entry. However, it may be challenging to distinguish new initiators from prevalent users of a drug in EHR data because some patients may have inadequate or no baseline data to determine prior drug exposure, especially for inpatient treatment studies. In other words, misclassification of prior use can occur if such use is recorded in other EHR systems for patients cared for by providers using different EHR systems. Some EHR systems have medication reconciliation data that are routinely recorded at specific medical encounters, including office visits, on hospital admission, and at discharge from hospitalization, where the providers record the medications that patients take at home, including the ones not prescribed or dispensed from the EHR system. This additional medication information, when accurate, may enhance the identification of prior drug exposure and reduce misclassification of prevalent vs. new users. Furthermore, it is worth noting that the definition of new users may considerably vary between studies given variation in data availability, and that each drug may have its unique time-varying hazard function, metabolism, and clearance profiles. The specific lengths of the washout period to define new initiators (i.e., new use is defined as no use in this washout period prior to the index use) should be based on pharmacodynamics and pharmacokinetics, and historical data availability, as well as the prescription pattern of each study drug.

- **The complexity of drug repurposing for COVID-19:** Several drugs assessed for utility in the prevention or treatment of COVID-19 were originally indicated for other conditions but repurposed for COVID-19 (i.e., “off-label use”).³⁴ Drug repurposing brings additional challenges to balancing the confounders at baseline because the same drug may be used to treat COVID-19 or the original indications (e.g., some may use angiotensin-converting enzyme inhibitors to prevent COVID-19 infection, whereas others use it to treat hypertension). Another challenge that may affect the specificity of exposure definitions is the frequent switching and discontinuation of therapies during the earlier phase in the pandemic when evidence was sparse and guidelines were rapidly evolving. One strategy is to use time-varying exposure definitions, such as dynamic treatment strategies with proper adjustment for time-varying confounding by marginal structural models or other g-methods.³⁵ It is also important to account for the reasons (usually captured as unstructured data) that give rise to the switch or discontinuation in these models.^{36,37}

Defining outcomes relevant for COVID-19

- **The opportunity of collecting data on specific clinical outcomes:** Although COVID-19 primarily affects the lungs, causing interstitial pneumonitis and severe acute respiratory distress syndrome (ARDS), it also affects multiple organs. A growing literature has identified some of the short- and long-term effects of COVID-19 on key markers of dysfunction in several organ systems (respiratory, cardiovascular, immune, musculoskeletal, hepatic, renal, and neurological). The rich and comprehensive clinical data contained in the EHRs often grant opportunities to ascertain these outcomes not only based on diagnosis and procedure codes but also abnormal vital signs, laboratory test results, or imaging findings.³⁸ In addition, for newly created diagnosis and procedure codes, such as ICD diagnosis code of COVID-19 and its complications, it is also possible to use EHR data for validation of the outcome definitions by chart review.
- **The challenge of standardizing clinical endpoints across EHR systems:** Many COVID-19 CER studies investigated hospitalization, intubation, intensive care unit (ICU) admission, and death as the outcomes of interest. These events are typically well-captured in EHR data sources, except for out-of-hospital death data, for which linkage to death records is often recommended. Attention should be paid to the interpretation of these outcomes as proxies of COVID-19 disease progression or recovery in the real-world setting, in the rapidly progressing pandemic, as some of these outcomes, such as ICU admission and oxygen use, might not be routinely collected in all data sources and influenced by the institute care protocol, hospital capacity, or supply. For example, ICU admission can be misclassified due to some units being repurposed as an ICU in response to the patient surge and potentially driven by the incidence rates in each period, in which case specific interventions that indicate critical illness may be a more reliable outcome, such as mechanical ventilation, extracorporeal membrane oxygenation, or use of vasopressors.

- **The lack of harmonization in the collection and reporting of outcomes:** In EHRs, the lack of direct comparability of results from studies across different health systems and geographic areas due to lack of harmonization in the data collection and reporting is a well-documented challenge and this is not unique in COVID-19 research. However, in situations like the COVID-19 pandemic, comparability in outcome collection and definition is a critical issue that must be addressed to halt the unprecedented havoc on public health and economies. In response to this issue, bodies such as the World Health Organization (WHO) have produced guidance on the minimum set of common outcome measures for studies of COVID-19 with the aim of enabling direct comparability and replication of CER studies across different settings.³⁹ Several consortiums have also been established to propose approaches for the aggregation of EHRs and related data to answer COVID-19 research questions. Some of these approaches suggested introducing specific diagnosis and procedure codes for identifying the COVID-19-related outcomes and establishing robust cohort definitions to ensure reproducibility and harmonization of concepts across different care settings.⁴⁰ For instance, the core outcomes developed by the WHO research group include viral burden (quantitative PCR or cycle threshold), patient survival (mortality at hospital discharge or at 60 days), and disease progress (hospital stay length and need for mechanical ventilation) and focus on the acute phase of COVID-19, whereas routinely collected safety data in EHRs, such as QT-prolongation and diagnosis of arrhythmias, may provide additional information regarding the CER of treatments in diverse patient populations.⁴¹ Currently, there is a lack of consensus (standard framework) on how to best evaluate long-term effects of COVID-19, including ambiguity in defining “long COVID” related and concurrent disorders and whether these are related to the disease itself or as a result of therapeutic/vaccines safety outcomes.⁴² Ongoing research funding has been made available to support research into “long COVID,” including developing an EHR-based registry detailing symptoms linked to patients’ samples that may further characterize this disease.⁴³

Minimizing common biases

- **Confounding bias:** As previously noted, EHRs contain rich clinical data typically not available in insurance claims data; many of which are potential confounders in COVID-19 CER, including vital signs (e.g., oxygen saturation, blood pressures, and body temperature), lifestyle factors (e.g., body mass index, smoking status, and alcohol consumption⁴⁴), laboratory tests (e.g., C-reactive protein and lactate dehydrogenase, D-dimer, and other biomarkers for inflammation or disease severity), and imaging findings (e.g., chest X-ray or computed tomography evidence of pulmonary infiltration or embolism).⁴⁵ In a rapidly evolving pandemic like COVID-19, it is crucial to adjust for potential confounding by calendar time trends in clinical practices. It is also important to consider changes in data availability and quality over time.⁴⁶ Besides, much of the essential confounder information, such as patient-reported symptoms, severity, stage, prognosis of disease, and functional status,⁴⁷ is recorded in

free-text notes or reports in EHRs, although this may not be consistent across hospital and EHR systems. Although substantially underutilized for confounding adjustment, adding unstructured information can potentially enhance researchers' ability to reduce confounding after using NLP to convert the free-text data into an analyzable format.^{18,37} Differences in interventions between health settings could also be explored using advanced techniques, such as the use of high-dimensional propensity scores with machine learning algorithms or instrumental variable analysis, can be considered to adjust for proxies of unmeasured confounding.⁴⁸ However, for instrumental variable analyses, it is challenging to identify a valid instrument and often requires strong assumptions. For instance, prescriber preference has been used as a potential instrumental variable, but if the preference of different prescribers is linked to their quality of care that is associated with the outcome of interest (which is often the case), the assumption of the instrumental variable being only linked to the outcome through the treatment is violated.⁴⁹ Furthermore, advanced machine-learning algorithms can be helpful to model the confounders either using propensity scores or treatment effects on the outcomes but in the context of COVID-19 research, sometimes we are limited by the number of users (e.g., newly available treatments/outcomes). In these cases, researchers should consider other techniques (dimension reduction and oversampling technique).⁵⁰

- **Missing data:** To properly handle missing data, investigators need to understand the mechanism of data missingness. Missing data may occur (1) “missing completely at random” (i.e., missingness is independent of all factors; e.g., missing a batch of laboratory results due to fire or a natural disaster); (2) “missing at random” (i.e., missingness is only dependent on observed data; e.g., missing laboratory results in the rehabilitation facilities but no other facilities when the type of facilities is observed); (3) “missing not at random” (i.e., missingness is dependent on unobserved data; e.g., missing a specific laboratory test and or imaging results due to differential ordering pattern of the physicians but the reasons underlying the decisions are not measured). Under missing completely at random, performing analysis using only those with complete data will not result in bias but may reduce statistical power. Under missing at random, investigators need to collect and adjust for these factors underlying missingness using proper methods (e.g., multiple imputation, maximum likelihood-based methods, or inverse probability weighting).⁵¹ However, these methods are less appropriate when the prevalence of missing data is very high. Under missing not at random, bias is generally expected, and investigators should attempt to assess the magnitude of such impact on the study estimates.⁵²

In studies using RWD, investigators typically assume the absence of recording of a disease state (e.g., having a negative test for COVID-19 diagnosis) as the absence of the condition, thus EHR-based CER often turns missing data into misclassification of the study variables. For example, previous studies in COVID-19 research showed a significant proportion of missing data on pre-existing health conditions to allow calculation of comorbidities

for the included patients.^{53,54} Unlike claims data in which the enrollment of the insurance coverage has well-documented start and end-dates, there is no “enrollment” or “membership” defined in an EHR. The EHR discontinuity (e.g., receiving care outside of a particular EHR system) has been shown to be associated with a large amount of information bias in essential variables in CER.¹⁷ Applying a prediction model to identify patients with high EHR-continuity and restrict the analysis among these patients can substantially reduce such biases.⁵⁵ A legitimate concern of this approach is whether findings based on those with high EHR-continuity are generalizable to the general population. To address this concern, a prior study has demonstrated that the patients with high EHR-continuity have similar comorbidity profiles compared to those with low EHR-continuity based on claims data that are not affected by EHR discontinuity.¹³ Data linkage of EHRs with other data sources (e.g., claims data with a shorter time lag, such as local insurance plan data or state-reported Medicaid data or leveraging novel electronic data collection methods, such as software application on smartphones, to capture data in the real-world setting) is important to address information bias due to EHR data discontinuity, although this process is often complicated by privacy concerns (e.g., the requirement of patient identifiers for data linkage), different clinical terminologies, technical specifications, and functional capabilities of different data sources.^{56,57}

- **Selection bias (or collider bias):** It can occur if restricting an analysis to those who had a cohort-qualifying event, such as hospitalization with COVID-19, had been tested for active infection or who have volunteered their participation in a prospective study (i.e., conditioning on a collider variable).⁵⁸ It can also happen with studies that included only patients without missing data when the missingness did not occur completely at random (i.e., “no missing data” effectively becomes the cohort inclusion criterion). The spurious association is expected if the collider variable is simultaneously associated with the treatment and outcome of interest. Such bias can be addressed to a certain extent by inverse probability weighting with the weights being the reciprocal of the probability of being selected into the cohort, conditioning on the predictors of the cohort-qualifying event.³⁶

The issues noted above are specific concerns that threaten the internal validity of studies based on data from EHRs. However, it is worth noting that, in COVID-19 research, there are specific challenges linked to generalizing the findings from EHR-based studies beyond the sample in these studies. Some of these challenges are related to geographical restrictions by catchment area or hospital setting. Researchers should explicitly define unique network features of databases (e.g., academic, outpatient-only, and hospital-based), such as selective entrance into the healthcare system (e.g., severity of COVID-19 disease), variations in care (e.g., admission criteria to ICU), and interpret the findings of these studies for similar study populations (in terms of disease status and other patient characteristics). Applications of external validity, particularly in the context of CER estimates, are only justified when the sample is representative of the population to which results are to be

generalized and, for that reason, nationwide cohorts enhance the generalizability and precision of findings and allow a wide range of research studies to be conducted. In the data analysis phase, we also encourage researchers to use other external data sources that can help to quantify the bias(es) arising from differences in the sample and the target population.

OPERATIONAL CHALLENGES USING EHRs IN COVID-19 CER

Operational challenges due to the COVID-19 pandemic, differences in the healthcare systems response, and in data recording across hospitals may hinder the comparability of outcomes across different settings. These challenges may further obstruct the possibility of combining data from different databases. Geographic variation in patient management strategies and the inability to capture exposure information from EHR data sources can also be an important challenge. Different countries or regions within countries may adopt different strategies regarding the initiation of treatment in inpatient or outpatient settings. In addition, the availability of EHRs that are more readily available for research differs from country to country. A previous study has shown the importance of differentiating between categories of patients admitted to hospitals and triaged to home; these care choices may not reflect similar patient physiology but instead reflect local care provision.⁵⁹ Data measurement (e.g., safety outcomes) or detailed record-keeping on patients' regular monitoring may also significantly be impacted by the availability of medical staff and the emergency caused by the unpredicted number of patients admitted with COVID-19 across different settings. During patient surge in the COVID-19 pandemic, many hospitals are understaffed. The increased patient-to-staff ratios could lead to a reduced quality in clinical care but also in data documentation and information recording. For instance, vital signs, height, weight, smoking status, may not get recorded as regularly or as reliably in a pandemic, which may lead to information bias or exacerbate missing data. Therefore, investigators should conduct more detailed evaluation of quality metrics of the EHR data generated during the pandemic, calibrated by the pre-pandemic data when possible. If major quality issues are identified for certain EHR components that are essential for study validity, the research team should consider linkage between EHR with additional data sources.³⁸ In addition, it is possible that some medical care of the study participants was provided in another EHR system and not captured by the study or that variation exists in how different sites extracted and recorded data in their systems (for example, differences in data recording between clinical notes and electronic systems especially at the beginning of the pandemic). Filling the data gaps by generating linkable identifiers is critical to address mismeasurement and discontinuity of care provision due to the lack of a centralized healthcare system between healthcare providers (hospitals, nursing homes, and general practitioners). However, caution should be paid around the potential risk of selection bias caused by incomplete linkage.

Last, it is also important to assess treatment effect heterogeneity by patient characteristics, care setting, and time trend, considering changes in clinical practice over time and with variants of concerns. There is abundant literature on EHR-based prognostic prediction

models, which can be informative to define different risk groups based on variables available in an EHR.⁶⁰ Given the substantial differences in public health policy, care delivery systems, and EHR data structures, when combining results from analyses conducted in several EHRs or health system data platforms, it is often recommended to stratify CER analyses by geography, healthcare systems, and study databases. The expected large underlying effect of heterogeneity can also be addressed by adopting random effect models for meta-analyses. Recently, several extensions to these models have been proposed to allow for heterogeneity across methods of imputation and adjustment for measurement errors in COVID-19 research.⁶¹

SUMMARY AND CONCLUSIONS

The COVID-19 pandemic has presented an unprecedented need for timely and reliable assessment of safety and effectiveness of therapeutic and preventive interventions and the wide availability of RWD can play a significant role in the generation of useful knowledge.

- EHRs represent one important source of RWD that may be critical in developing RWE to inform healthcare decision making for COVID-19 without the need for primary data collection, something that would further negatively impact an already overburdened healthcare system.
- Whereas common principles to avoid information bias, selection bias, confounding bias, and model misspecification remain applicable to EHR-based CER, the rapidly evolving pandemic hitting an overburdened healthcare system could exaggerate these common biases in nonrandomized studies. Therefore, the study findings should be interpreted cautiously, and study limitations should be acknowledged explicitly. Engagement with systems generating the data will provide an important insight of the data origins and data gaps.
- The research team should start with careful construction of the research questions which requires setting up clear definitions of the study population based on treatment exposure and proper choice of the comparator groups, considering the evolving knowledge about COVID-19, and understanding the rapidly changing clinical practice patterns over time.
- EHRs contain rich clinical data for assessing relevant end points, and the identification of potential key confounders and effect modifiers relevant for answering the specific CER questions regarding therapeutic and vaccine interventions for COVID-19. However, proper data processing (such as transforming free-text unstructured EHR data into an analyzable data set) and quality checking (such as assessing impact of data discontinuity and missing data) are often warranted.
- Pooled analyses across multiple EHR systems are often needed to accrue sufficient power and to demonstrate the generalizability of the study findings across settings. Data harmonization and outcome reporting standardization across sites are pivotal to ensure study validity.

In conclusion, EHRs provide an opportunity to perform rapid COVID-19 CER due to availability of both structured and

unstructured data, such as laboratory and imaging data, as well as relevant confounders, independent risk factors, and potential for data linkages to create a holistic view of patient management and outcomes. However, unique features of COVID-19, including varying disease presentation, disease measurement, and constantly changing clinical management during an emerging pandemic, require special considerations to produce reliable evidence to support healthcare decision making. Given the potential challenges to study validity, we should interpret the nonrandomized RWE considering ongoing RCTs. With proper design and analysis, EHR-based CER can be helpful for (1) RCT hypothesis generation: identifying existing pharmacotherapies that could potentially be repurposed to treat COVID-19, and (2) real-world effectiveness determination: assessing the causal treatment effects of COVID-19 treatments and vaccines in the vulnerable populations substantially underrepresented by the RCTs, such as patients with advanced kidney and liver diseases, or those who are severely immunocompromised or frail.^{62,63} The COVID-19 pandemic highlights not only the challenges in using EHRs to inform prescribing decisions, but also its unique potential to generate generalizable information from real-world heterogeneous populations. This paper provided structured guidance for the proper conduct and appraisal of drug and vaccines' effectiveness and safety research using EHRs for the ongoing and future pandemics.

SUPPORTING INFORMATION

Supplementary information accompanies this paper on the *Clinical Pharmacology & Therapeutics* website (www.cpt-journal.com).

ACKNOWLEDGMENTS

This paper is endorsed by the International Society for Pharmacoepidemiology (ISPE). We acknowledge the support of research assistants (Winnie Ho and Luke Edward Zobotka) in searching and reviewing some of the supporting materials for the development of this manuscript. We also thank Dr. Gianluca Trifiro, a member of the International Society for Pharmacoepidemiology, for his participation.

FUNDING

Funding for the development of this paper was provided by ISPE (<https://www.pharmacoepi.org/ISPE/assets/File/General/Manuscript%20Proposals%20Previously%20Accepted.pdf>).

CONFLICT OF INTEREST

G.S. is employed by Cytel Inc. DB is an employee of Takeda. A.R.Z. has received grant funding from Sanofi Pasteur direct to Brown University for research on the epidemiology of infections and vaccinations among nursing home residents and infants. A.D.L. is an employee of Boehringer Ingelheim International. M.S.G. is an employee of Bayer AG and owns stock in the company. Y.Y. is an employee of AbbVie and may own AbbVie stock or options. J.L.'s spouse is an employee of GlaxoSmithKline and owns stock in the company. J.A.L. is an employee of IQVIA and owns stock in the company. M.G. is an employee of Merck and owns stock in the company. H.Y. has research grant from Pfizer for unrelated work. T.D. provides consulting services to pharma industry. C.T.R., X.L., W.L., X.W., C.Y., J.K.L., and D.C.M. declared no competing interests for this work.

DISCLAIMER

This paper reflects the views and opinions of the authors and inputs that members of the working group shared in their personal capacity and do not represent the position of their respective bodies. The opinions expressed in this manuscript are those of the authors and should not be interpreted as the position of the US Food and Drug Administration.

COLLABORATORS

Guidance was received from the Working Group of the International Society for Pharmacoepidemiology Comparative Effectiveness Research Special Interest Group.

© 2022 The Authors. *Clinical Pharmacology & Therapeutics* published by Wiley Periodicals LLC on behalf of American Society for Clinical Pharmacology and Therapeutics.

This is an open access article under the terms of the Creative Commons Attribution-NonCommercial-NoDerivs License, which permits use and distribution in any medium, provided the original work is properly cited, the use is non-commercial and no modifications or adaptations are made.

1. WHO announces COVID-19 outbreak a pandemic [Internet]. <https://www.euro.who.int/en/health-topics/health-emergencies/coronavirus-covid-19/news/news/2020/3/who-announces-covid-19-outbreak-a-pandemic>. Accessed August 26, 2021.
2. WHO Coronavirus (COVID-19) Dashboard [Internet]. <https://covid19.who.int>. Accessed October 10, 2021.
3. Fegert, J.M., Vitiello, B., Plener, P.L. & Clemens, V. Challenges and burden of the Coronavirus 2019 (COVID-19) pandemic for child and adolescent mental health: a narrative review to highlight clinical and research needs in the acute phase and the long return to normality. *Child Adolescent Psychiatry Mental Health* **14**, 20 (2020).
4. Commissioner. FDA Approves First COVID-19 Vaccine [Internet]. FDA. FDA; 2021. <https://www.fda.gov/news-events/press-announcements/fda-approves-first-covid-19-vaccine>. Accessed October 10, 2021.
5. Regulatory approval of Pfizer/BioNTech vaccine for COVID-19 [Internet]. GOV.UK. <https://www.gov.uk/government/publications/regulatory-approval-of-pfizer-biontech-vaccine-for-covid-19>. Accessed October 10, 2021.
6. PINHO AC. EMA recommends first COVID-19 vaccine for authorisation in the EU [Internet]. European Medicines Agency. 2020. <https://www.ema.europa.eu/en/news/ema-recommends-first-covid-19-vaccine-authorisation-eu>. Accessed October 10, 2021.
7. Chakraborty, I. & Maity, P. COVID-19 outbreak: Migration, effects on society, global environment and prevention. *Sci. Total Environ.* **728**, 138882 (2020).
8. Dean, N.E. et al. Creating a framework for conducting randomized clinical trials during disease outbreaks. *N. Engl. J. Med.* **382**, 1366–1369 (2020).
9. Radawski, C.A. et al. The utility of real-world evidence for benefit-risk assessment, communication, and evaluation of pharmaceuticals: Case studies. *Pharmacoepidemiol. Drug Saf.* **29**, 1532–1539 (2020).
10. Kruse, C.S., Stein, A., Thomas, H. & Kaur, H. The use of electronic health records to support population health: a systematic review of the literature. *J. Med. Syst.* **42**, 214 (2018).
11. Franklin, J.M., Lin, K.J., Gatto, N.M., Rassen, J.A., Glynn, R.J. & Schneeweiss, S. Real-world evidence for assessing pharmaceutical treatments in the context of COVID-19. *Clin. Pharmacol. Ther.* **109**, 816–828 (2021).
12. Brat, G.A. et al. International electronic health record-derived COVID-19 clinical course profiles: the 4CE consortium. *NPJ Digit. Med.* **3**, 1–9 (2020). <https://www.nature.com/articles/s41746-020-00308-0>
13. Lin, K.J. et al. External validation of an algorithm to identify patients with high data-completeness in electronic health records for comparative effectiveness research. *Clin. Epidemiol.* **12**, 133–141 (2020).
14. Madhavan, S. et al. Use of electronic health records to support a public health response to the COVID-19 pandemic in the United States: a perspective from 15 academic medical centers. *J. Am. Med. Inform. Assoc.* **28**, 393–401 (2021).
15. Two elite medical journals retract coronavirus papers over data integrity questions | Science | AAAS [Internet]. <https://www.sciencemag.org/news/2021/09/01/two-elite-medical-journals-retract-coronavirus-papers-over-data-integrity-questions>

- cemag.org/news/2020/06/two-elite-medical-journals-retract-coronavirus-papers-over-data-integrity-questions. Accessed August 26, 2021.
16. Wells, B.J., Chagin, K.M., Nowacki, A.S. & Kattan, M.W. Strategies for handling missing data in electronic health record derived data. *EGEMS (Wash DC)* **1**, 1035 (2013).
 17. Lin, K.J., Glynn, R.J., Singer, D.E., Murphy, S.N., Lii, J. & Schneeweiss, S. Out-of-system care and recording of patient characteristics critical for comparative effectiveness research. *Epidemiology* **29**, 356–363 (2018).
 18. Luo, Y. et al. Natural language processing for EHR-based pharmacovigilance: a structured review. *Drug Saf.* **40**, 1075–1089 (2017).
 19. Burn, E., et al. Deep phenotyping of 34,128 patients hospitalised with COVID-19 and a comparison with 81,596 influenza patients in America, Europe and Asia: an international network study. *medRxiv Preprint*. June 28, 2020; <https://doi.org/2020.04.22.20074336>. [e-pub ahead of print].
 20. Dinnes, J. et al. Rapid, point-of-care antigen and molecular-based tests for diagnosis of SARS-CoV-2 infection. *Cochrane Database Syst. Rev.* **3**, CD013705 (2021).
 21. Scourfield, D.O. et al. The role and uses of antibodies in COVID-19 infections: a living review. *Oxf. Open Immunol.* **2**, iqab003 (2021).
 22. Weiskopf, N.G., Hripsak, G., Swaminathan, S. & Weng, C. Defining and measuring completeness of electronic health records for secondary use. *J. Biomed. Inform.* **46**, 830–836 (2013).
 23. Baum, U., Kulathinal, S. & Auranen, K. Exposure misclassification bias in the estimation of vaccine effectiveness. *PLoS One* **16**, e0251622 (2021).
 24. Feng, S., Hategeka, C. & Grépin, K.A. Addressing missing values in routine health information system data: an evaluation of imputation methods using data from the Democratic Republic of the Congo during the COVID-19 pandemic. *Popul. Health Metr.* **19**, 44 (2021).
 25. Koné, S., Bonfoh, B., Dao, D., Koné, I. & Fink, G. Heckman-type selection models to obtain unbiased estimates with missing measures outcome: theoretical considerations and an application to missing birth weight data. *BMC Med. Res. Methodol.* **19**, 231 (2019).
 26. D'Arcy, M., Stürmer, T. & Lund, J.L. The importance and implications of comparator selection in pharmacoepidemiologic research. *Curr. Epidemiol. Rep.* **5**, 272–283 (2018).
 27. Rossi, B. et al. Effect of tocilizumab in hospitalized patients with severe COVID-19 pneumonia: a case-control cohort study. *Pharmaceuticals (Basel)* **13**, E317 (2020).
 28. Osborne, T.F., Veigulis, Z.P., Arreola, D.M., Mahajan, S.M., Röösl, E. & Curtin, C.M. Association of mortality and aspirin prescription for COVID-19 patients at the Veterans Health Administration. *PLoS One* **16**, e0246825 (2021).
 29. Lund, J.L., Richardson, D.B. & Stürmer, T. The active comparator, new user study design in pharmacoepidemiology: historical foundations and contemporary application. *Curr. Epidemiol. Rep.* **2**, 221–228 (2015).
 30. Anti-SARS-CoV-2 Monoclonal Antibodies | COVID-19 Treatment Guidelines [Internet]. <https://www.covid19treatmentguidelines.nih.gov/therapies/anti-sars-cov-2-antibody-products/anti-sars-cov-2-monoclonal-antibodies/>. Accessed August 26, 2021.
 31. Corticosteroids [Internet]. COVID-19 Treatment Guidelines. <https://www.covid19treatmentguidelines.nih.gov/therapies/immunomodulators/corticosteroids/>. Accessed August 26, 2021.
 32. Stürmer, T., Wang, T., Golightly, Y.M., Keil, A., Lund, J.L. & Jonsson Funk, M. Methodological considerations when analysing and interpreting real-world data. *Rheumatology (Oxford)* **59**, 14–25 (2020).
 33. Jung, R.G. et al. Methodological quality of COVID-19 clinical research. *Nat. Commun.* **12**, 943 (2021).
 34. Rentsch, C.T. et al. Early initiation of prophylactic anticoagulation for prevention of COVID-19 mortality: a nationwide cohort study of hospitalized patients in the United States. *medRxiv*. <https://doi.org/2020.12.09.20246579>.
 35. Li, X., Young, J.G. & Toh, S. Estimating effects of dynamic treatment strategies in pharmacoepidemiologic studies with time-varying confounding: a primer. *Curr. Epidemiol. Rep.* **4**, 288–297 (2017).
 36. Robins, J.M., Hernán, M.A. & Brumback, B. Marginal structural models and causal inference in epidemiology. *Epidemiology* **11**, 550–560 (2000).
 37. Naimi, A.I., Cole, S.R. & Kennedy, E.H. An introduction to g methods. *Int. J. Epidemiol.* **46**, 756–762 (2017).
 38. Dagliati, A., Malovini, A., Tibollo, V. & Bellazzi, R. Health informatics and EHR to support clinical research in the COVID-19 pandemic: an overview. *Brief. Bioinform.* **22**, 812–822 (2021).
 39. A minimal common outcome measure set for COVID-19 clinical research - PubMed [Internet]. <https://pubmed.ncbi.nlm.nih.gov/32539990/>. Accessed August 25, 2021.
 40. The COVID-19 Core Outcome Set Project | Cochrane Methods [Internet]. <https://methods.cochrane.org/news/covid-19-core-outcome-set-project-invitation-take-part>. Accessed August 26, 2021.
 41. WHO Working Group on the Clinical Characterisation and Management of COVID-19 infection. A minimal common outcome measure set for COVID-19 clinical research. *Lancet Infect. Dis.* **20**(8), e192–e197 (2020).
 42. Rando, H.M. et al. Challenges in defining Long COVID: Striking differences across literature, Electronic Health Records, and patient-reported information. *medRxiv Preprint*. <https://doi.org/2021.03.20.21253896>.
 43. NIH launches new initiative to study “Long COVID” [Internet]. National Institutes of Health (NIH). 2021. <https://www.nih.gov/about-nih/who-we-are/nih-director/statements/nih-launches-new-initiative-study-long-covid>. Accessed September 6, 2021.
 44. Umnuaypornlert, A., Kanchanasurakit, S., Lucero-Priso, D.E.I. & Saokaew, S. Smoking and risk of negative outcomes among COVID-19 patients: A systematic review and meta-analysis. *Tob. Induc. Dis.* **19**, 9 (2021).
 45. Gong, K. et al. A multi-center study of COVID-19 patient prognosis using deep learning-based CT image analysis and electronic health records. *Eur. J. Radiol.* **139**, 109583 (2021).
 46. Bayley, K.B., Belnap, T., Savitz, L., Masica, A.L., Shah, N. & Fleming, N.S. Challenges in using electronic health record data for CER: experience of 4 learning organizations and solutions applied. *Med. Care* **51**(8 Suppl 3), S80–S86 (2013).
 47. Katzenschlager, S. et al. Can we predict the severe course of COVID-19 – a systematic review and meta-analysis of indicators of clinical outcome? *PLoS One* **16**, e0255154 (2021).
 48. Zhang, X., Faries, D.E., Li, H., Stamey, J.D. & Imbens, G.W. Addressing unmeasured confounding in comparative observational research. *Pharmacoepidemiol. Drug Saf.* **27**, 373–382 (2018).
 49. Franklin, J.M., Schneeweiss, S., Huybrechts, K.F. & Glynn, R.J. Evaluating possible confounding by prescriber in comparative effectiveness research. *Epidemiology* **26**(2), 238–241 (2015). <https://doi.org/10.1097/ede.0000000000000241>
 50. Elreedy, D. & Atiya, A.F. A comprehensive analysis of synthetic minority oversampling technique (SMOTE) for handling class imbalance. *Inf. Sci.* **1**(505), 32–64 (2019).
 51. Carpenter, J.R. & Smuk, M. Missing data: A statistical framework for practice. *Biom J.* **63**, 915–947 (2021).
 52. Madden, J.M., Lakoma, M.D., Rusinak, D., Lu, C.Y. & Soumerai, S.B. Missing clinical and behavioral health data in a large electronic health record (EHR) system. *J. Am. Med. Inform. Assoc.* **23**, 1143–1149 (2016).
 53. DeLozier, S. et al. Phenotyping coronavirus disease 2019 during a global health pandemic: Lessons learned from the characterization of an early cohort. *J. Biomed. Inform.* **117**, 103777 (2021).
 54. Brown, J.S., Bastarache, L. & Weiner, M.G. Aggregating electronic health record data for COVID-19 research—caveat emptor. *JAMA Network Open* **4**, e2117175 (2021).
 55. Lin, K.J., Singer, D.E., Glynn, R.J., Murphy, S.N., Lii, J. & Schneeweiss, S. Identifying patients with high data completeness to improve validity of comparative effectiveness research in electronic health records data. *Clin. Pharmacol. Ther.* **103**, 899–905 (2018).

56. Lo, W.-C. *et al.* Enhancing data linkage to break the chain of COVID-19 spread: The Taiwan experience. *J. Med. Int. Res.* **23**, e24294 (2021).
57. Bhattacharya, A. *et al.* Healthcare-associated COVID-19 in England: a national data linkage study. *J. Infect.* **83**(5), 565–572 (2021). <https://doi.org/10.1016/j.jinf.2021.08.039>
58. Griffith, G.J. *et al.* Collider bias undermines our understanding of COVID-19 disease risk and severity. *Nat. Commun.* **11**, 5749 (2020).
59. Turcato, G., Zaboli, A. & Pfeifer, N. The COVID-19 epidemic and reorganisation of triage, an observational study. *Intern. Emerg. Med.* **9**, 1–8 (2020).
60. Wynants, L. *et al.* Prediction models for diagnosis and prognosis of covid-19: systematic review and critical appraisal. *BMJ* **369**, m1328 (2020)
61. Evaluation of COVID-19 vaccine effectiveness [Internet]. https://www.who.int/publications/i/item/WHO-2019-nCoV-vaccine_effectiveness-measurement-2021.1. Accessed September 14, 2021.
62. Pottegård, A., Kurz, X., Moore, N., Christiansen, C.F. & Klungel, O. Considerations for pharmacoepidemiological analyses in the SARS-CoV-2 pandemic. *Pharmacoepidemiol. Drug Saf.* **29**, 825–831 (2020).
63. Arlett, P., Kjær, J., Broich, K. & Cooke, E. Real-world evidence in EU medicines regulation: enabling use and establishing value. *Clin. Pharmacol. Ther.* **111**, 21–23 (2022).