

# International Electronic Health Record-Derived COVID-19 Clinical Course Profile: The 4CE Consortium

Gabriel A Brat\*; Griffin M Weber\*; Nils Gehlenborg; Paul Avillach; Nathan P Palmer; Luca Chiovato; James Cimino; Lemuel R Waitman; Gilbert S Omenn; Alberto Malovini; Jason H Moore; Brett K Beaulieu-Jones; Valentina Tibollo; Shawn N Murphy; Sehi L'Yi; Mark S Keller; Riccardo Bellazzi; David A Hanauer; Arnaud Serret-Larmande; Alba Gutierrez-Sacristan; John J Holmes; Douglas S Bell; Kenneth D Mandl; Robert W Follett; Jeffrey G Klann; Douglas A Murad; Luigia Scudeller; Mauro Bucalo; Katie Kirchoff; Jean Craig; Jihad Obeid; Vianney Jouhet; Romain Griffier; Sebastien Cossin; Bertrand Moal; Lav P Patel; Antonio Bellasi; Hans U Prokosch; Detlef Kraska; Piotr Sliz; Amelia LM Tan; Kee Yuan Ngiam; Alberto Zambelli; Danielle L Mowery; Emily Schiver; Batsal Devkota; Robert L Bradford; Mohamad Daniar; APHP/Universities/INSERM COVID-19 research collaboration; Christel Daniel; Vincent Benoit; Romain Bey; Nicolas Paris; Patricia Serre; Nina Orlova; Julien Dubiel; Martin Hilka; Anne Sophie Jannot; Stephane Breant; Judith Leblanc; Nicolas Griffon; Anita Burgun; Melodie Bernaux; Arnaud Sandrin; Elisa Salamanca; Thomas Ganslandt; Tobias Gradinger; Julien Champ; Martin Boeker; Patricia Martel; Alexandre Gramfort; Olivier Grisel; Damien Leprovost; Thomas Moreau; Gael Varoquaux; Jill Jen Vie; Demian Wassermann; Arthur Mensch; Charlotte Caucheteux; Christian Haverkamp; Guillaume Lemaitre; Christian Haverkamp; The Consortium for Clinical Characterization of COVID-19 by EHR (4CE); Tianxi Cai†; Isaac S Kohane†

\* authors contributed equally. † co-corresponding authors. ^ authors 6-83 ordered at random. Affiliations follow references.

## ABSTRACT

**INTRODUCTION:** The Coronavirus Disease 2019 (COVID-19) epidemic has caused extreme strains on health systems, public health infrastructure, and economies of many countries. A growing literature has identified key laboratory and clinical markers of pulmonary, cardiac, immune, coagulation, hepatic, and renal dysfunction that are associated with adverse outcomes. Our goal is to consolidate and leverage the largely untapped resource of clinical data from electronic health records of hospital systems in affected countries with the aim to better-define markers of organ injury to improve outcomes.

**METHODS:** A consortium of international hospital systems of different sizes utilizing Informatics for Integrating Biology and the Bedside (i2b2) and Observational Medical Outcomes Partnership (OMOP) platforms was convened to address the COVID-19 epidemic. Over a course of two weeks, the group initially focused on admission comorbidities and temporal changes in key laboratory test values during infection. After establishing a common data model, each site generated four data tables of aggregate data as comma-separated values files. These non-interlinked files encompassed, for COVID-19 patients, daily case counts; demographic breakdown; daily laboratory trajectories for 14 laboratory tests; and diagnoses by diagnosis codes.

**RESULTS:** 96 hospitals in the US, France, Italy, Germany, and Singapore contributed data to the consortium for a total of 27,927 COVID-19 cases and 187,802 performed laboratory values. Case counts and laboratory trajectories were concordant with existing literature. Laboratory test values at the time of viral diagnosis showed hospital-level differences that were equivalent to country-level variation across the consortium partners.

**CONCLUSIONS:** In under two weeks, we formed an international community of researchers to answer critical clinical and epidemiological questions around COVID-19. Harmonized data sets analyzed locally and shared as aggregate data has allowed for rapid analysis and visualization of regional differences and global commonalities. Despite the limitations of our datasets, we have established a framework to capture the trajectory of COVID-19 disease in various subsets of patients and in response to interventions.

## Introduction

The Coronavirus Disease 2019 (COVID-19) pandemic has caught the world off guard, reshaping ways of life, the economy, and healthcare delivery all over the globe. The virulence and transmissibility of responsible virus (SARS-CoV-2) is striking. Crucially, there remains a paucity of relevant clinical information to drive response at the clinical and population levels. Even in an information technology-dominated era, fundamental measurements to guide public health decision-making remain unclear. Knowledge still lags on incidence, prevalence, case-fatality rates, and clinical predictors of disease severity and outcomes. While some of the knowledge gaps relate to the need for further laboratory testing, data that should be widely available in electronic health records have not yet been effectively shared across clinical sites, with public health agencies, or with policy makers. At the time of this writing, more than three months after the earliest reports of the disease in China, only 5.8% of US cases reported to the CDC have clinical details included <sup>1</sup>.

Even before trials are implemented to determine which therapies will provide, frontline clinicians are not yet benefitting from knowledge as basic as understanding the differences in the clinical course between male and female patients <sup>2</sup>. Through case studies and series, we have learned that COVID-19 can have multi-organ involvement. A growing literature has identified key markers of cardiac,<sup>3</sup> immune,<sup>4</sup> coagulation,<sup>5</sup> muscle,<sup>5,6</sup> hepatic,<sup>7</sup> and renal<sup>8</sup> injury and dysfunction, including extensive evidence of myocarditis and cardiac injury associated with severe disease. Laboratory perturbations in lactate dehydrogenase (LDH), C-reactive protein (CRP), and procalcitonin<sup>9</sup> have been described. However, data from larger cohorts, linked to outcomes, remain unavailable.

Because electronic health records (EHRs) are not themselves agile analytic platforms, we have been successfully building upon the open source and free i2b2 (for Informatics for Integrating Biology and the Bedside) toolkit <sup>10-17</sup> to manage, compute, and share data extracted from EHRs. In response to COVID-19, we have organized a global community of researchers, most of whom are or have been members of the i2b2 Academic Users Group, to rapidly set up an *ad hoc* network that can begin to answer some of the clinical and epidemiological questions around COVID-19 through data harmonization, analytics, and visualizations. The Consortium for Clinical Characterization of COVID-19 by EHR (4CE)—pronounced “foresee”—comprises partner hospitals from five countries.

Our early efforts aim to consolidate, share, and interpret data about the clinical trajectories of the infection in patients with a first focus on laboratory values and comorbidities. This initial report seeks (a) to establish the accessibility and suitability of data from electronic medical record for COVID-19 patients; (b) to learn about the clinical trajectories of patients; (c) to facilitate evaluation and communication of the utility of various laboratory tests and therapies; and (d) to contribute data, reproducible data mining and visualization workflows, and learnings to a global network and the broader public.

Here, we report on initial results and the structure of a new, rapidly formed network designed to be a highly scalable system, now implemented at 21 sites. The international scope of our collaboration allows us to identify some of the similarities in clinical course and a few country-specific variations. We recognize that these early data are incomplete and are subject to many biases and limitations, which constrain the conclusions we can currently draw. However, we believe the sources of our data and the mechanism we have established for sharing

them are sound, reproducible, and scalable. We also hope our results to-date will encourage other sites to share data and contribute to this important research effort.

## Methods

### Selection of Laboratory Values

Multiple studies have reported significant abnormalities in several laboratory tests in patients with COVID-19. Studies have shown abnormalities in cardiac, hepatic, renal, immune, and coagulation physiology. Those laboratory results are associated with both disease presentation and severity of disease. For this initial study, we decided to select a subset of laboratories that are commonly performed, as identified by the Regenstrief Institute responsible for the Logical Objects, Identifiers, Names and Codes (LOINC) standard,<sup>18</sup> and had been previously associated with worse outcomes in COVID-19 positive patients. Based on the meta-analysis of Lippi and Plebani,<sup>19</sup> we chose to focus on 14 laboratory studies that are commonly performed: alanine aminotransferase (ALT), aspartate aminotransferase (AST), total bilirubin (Tbili), albumin, cardiac troponin (high sensitivity), lactate dehydrogenase (LDH), D-dimer, white blood cell count (WBC), lymphocyte count, neutrophil count, procalcitonin, and prothrombin time. LOINC codes were identified for each laboratory study as well as the units and reference ranges.

### Cohort Identification

All patients who received a polymerase chain reaction (PCR) confirmed diagnosis of COVID-19 were included in the data collection. Some hospitals only included patients who were admitted to the hospital while others included all patients for whom the test was positive.

### Data Collection and Aggregation

Sites obtained the data for their files in several ways. Most sites leveraged the open source i2b2 software platform already installed at their institution<sup>20</sup> which supports query and analysis of clinical and genomics data. More than 200 organizations worldwide use i2b2 for a variety of purposes, including identifying patients for clinical trials, drug safety monitoring, and epidemiology research. Most 4CE sites with i2b2 used database scripts to directly query their i2b2 repository to calculate counts needed for data files. Institutions without i2b2 used their own clinical data warehouse solutions and querying tools to create the files. In some cases, a hybrid method was used that leveraged different data warehouse platforms to fill in i2b2 gaps. For example, Assistance Publique – Hôpitaux de Paris (APHP), the largest hospital system in Europe, aggregates all EHR data from 39 hospitals in Paris and its surroundings. APHP exported data from the Observational Medical Outcomes Partnership (OMOP) Common Data Model for transformation to the shared format.

Each site generated four data tables, saved as comma-separated values (CSV) files. To protect patient privacy, the files we report contain only aggregate counts (no data on individual patients). In order to further protect patient identity, small counts were obfuscated (see below), since an aggregate count of “1” represents an individual patient. By computing these values locally and only sharing the aggregate data, sites were able to obtain institutional approval more rapidly.

The first file, *DailyCounts.csv*, contained one row per calendar date. Each row included the date, the number of new COVID-19 positive patients, the number of COVID-19 patients in an intensive care unit (ICU), and the number of new deaths from COVID-19.

The second file, *Demographics.csv*, contained counts of the total number of COVID-19 positive patients, broken down by gender and age group (0-2, 3-5, 6-11, 12-17, 18-25, 26-49, 50-69, 70-79, and 80+ years old).

The third file, *Labs.csv*, described the daily trajectories of select laboratory tests. Each row corresponded to a laboratory test (identified using a LOINC code) and the number of days since a patient had a positive COVID-19 test, ranging from -6 (one week before the test result) to 1 (the day of the test result) to N (the day the file was created). The values in each row are the number of patients who have a test result on that day and the mean and standard deviation of the test results.

The fourth file, *Diagnoses.csv*, lists all the diagnoses recorded in the EHR for COVID-19 positive patients, starting from one week before their positive COVID-19 test to the present, with the count of the number of patients with the corresponding ICD-9 or ICD-10 code.

Sites can optionally obfuscate the values in any of these files by replacing small counts with “-1”. Sites can indicate missing data or data that they are unable to obtain (e.g., whether patients are in an ICU) with “-2”.

Sites uploaded their files to a private shared folder. These files were then merged into four combined files that included the totals from the individual sites. Each value in the combined file has four components: (1) the number of sites with unmasked values; (2) the sum of those values; (3) the number of sites with obfuscated values; and (4) the sum of the obfuscation thresholds for those sites. For example, if five sites report values 25, 15, -1 (between 00 and 9 patients), -1 (between 00 and 4 patients), -1 (between 00 and 4 patients), then the combined file lists two unmasked sites with a total of 40 patients and three masked sites with up to  $9+4+4=17$  patients. From this, it can be inferred that there are between 40 and 57 patients. Given the large geographic distance between our sites, we assumed that each COVID-19 positive patient was only represented in one EHR. The combined *Labs.csv* file contains a weighted average (rather than the sum) of the unmasked mean test results from each site.

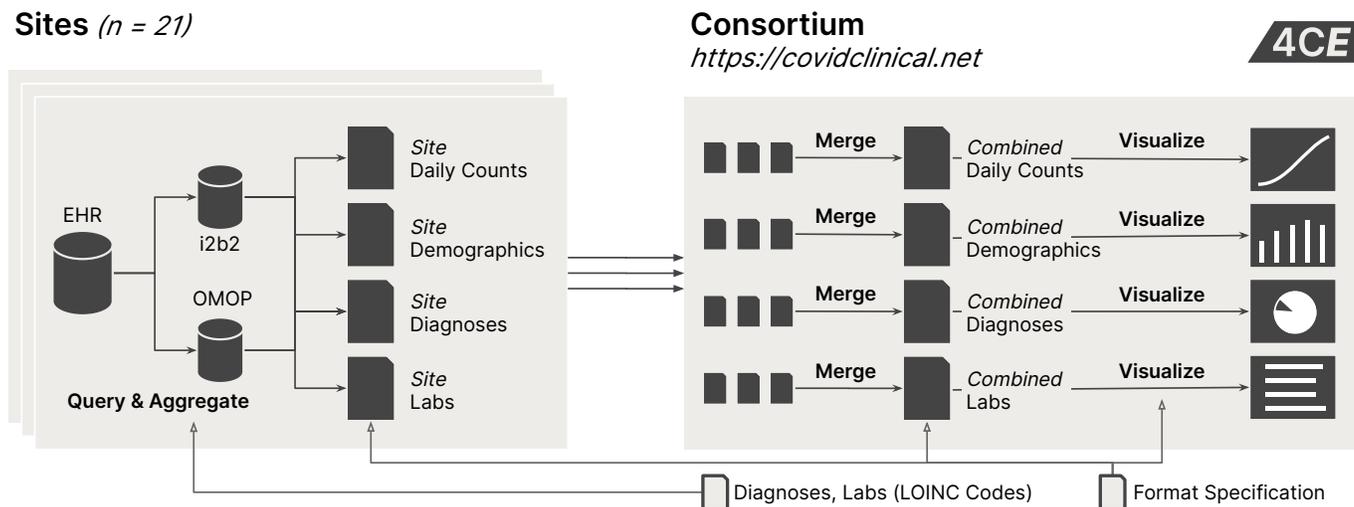
## ICD Mapping

Diagnosis codes were submitted from the sites as either international clinical diagnosis (ICD)-9 or ICD-10 billing codes. ICD-9 diagnosis codes were mapped to ICD-10 by first attempting to match the ICD-9 codes to child concepts of ICD-10 codes in the Accrual to Clinical Trials (ACT) ICD-10→ICD-9 ontology.<sup>21</sup> In the cases where no match was found in the ACT ontology, ICD-9 codes were matched to the ICD-10 codes that shared a common concept unique identifier (CUI) in the 2019 build of the US National Library of Medicine’s (NLM’s) Unified Medical Language System (UMLS).<sup>22</sup>

## Data Sharing and Visualization

We created a website hosted at <https://covidclinical.net> to provide interactive visualizations of our datasets as well as direct access to all shareable data collected for this publication. Data aggregation and publication

processes are shown in Figure 1. Visualizations were implemented using Python and Altair (<http://altair-viz.github.io/>) in Jupyter Notebooks (<https://jupyter.org/>), all of which are freely available on the website. The Vega visualizations (<http://vega.github.io>) generated by Altair were embedded into a Jekyll-based site (<http://jekyllrb.com/>) that is hosted on Amazon Web Services.



**Figure 1.** Overview of data collection and analysis.

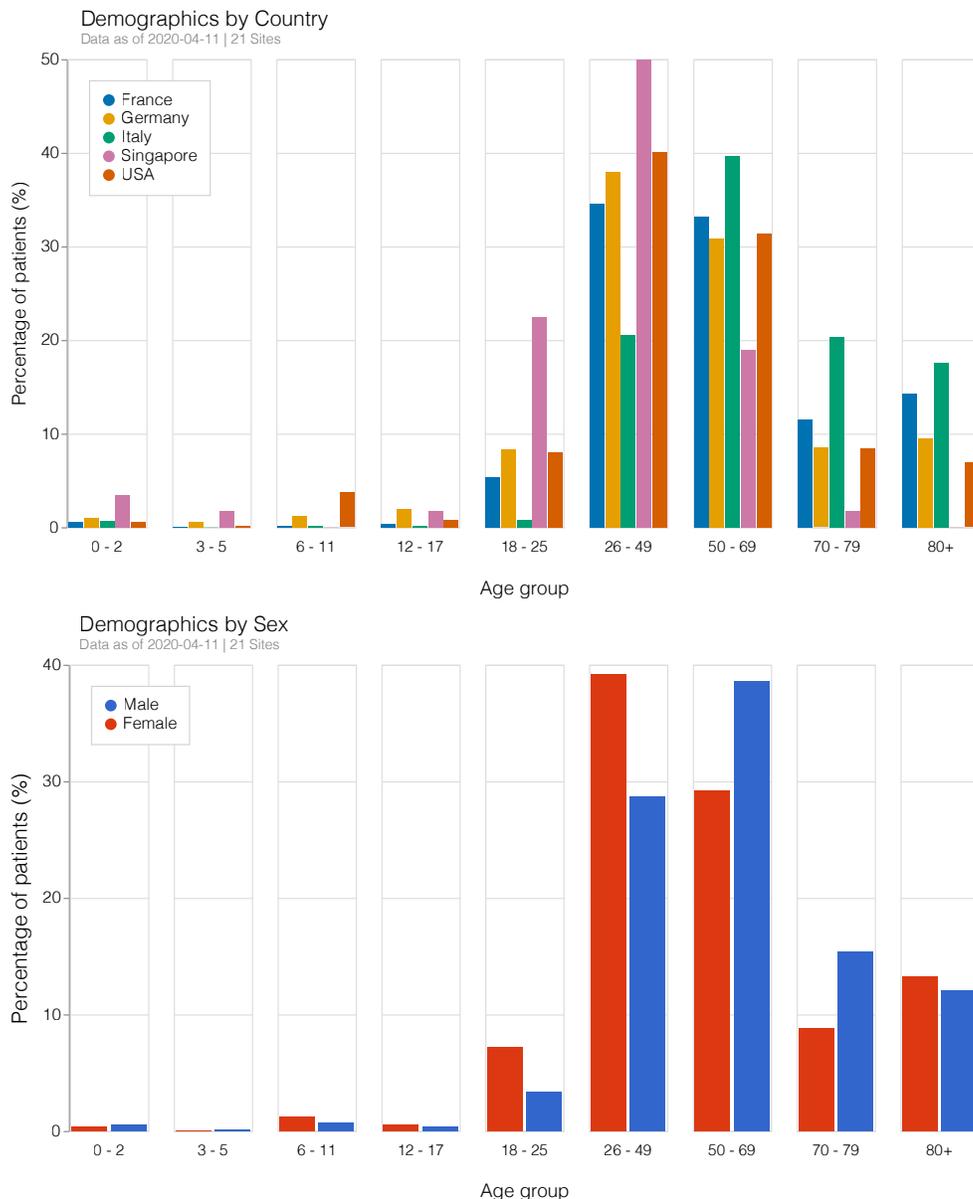
## Data Availability

Data files for daily counts, demographics, diagnosis, and labs datasets are available at <https://covidclinical.net>.

## Informed Consent/IRB Statement

Each institution reported obtaining proper institutional review board approval for data sharing. Certifications of waivers or approval were collected by the Consortium. As data were transmitted in aggregate, no patient level data were available from any site.

## Results



**Figure 2. (a) Patients by country and age group. (b) Patients by sex and age group.**

### Demographic and Consortium Level Data

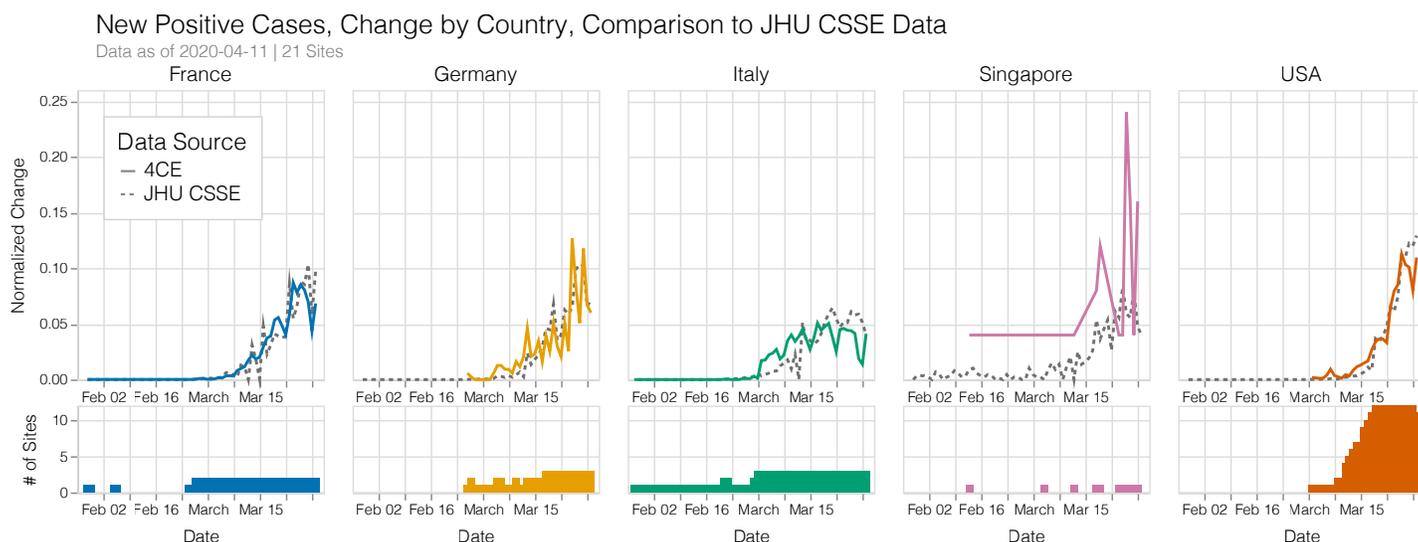
Over a span of three weeks, 96 total hospitals in the US (45), France (42), Italy (5), Germany (3), and Singapore (1) contributed data to the consortium. This was represented by 21 data collaboratives across these five countries. A total of 27,584 patients with COVID-19 diagnosis were included in the dataset, with data covering January 1, 2020 through April 9th, 2020. We collected 187,802 laboratory values and harmonized them across sites. 13.0% of sites submitted complete data sets that included values for each laboratory (39.1% for at least 13, and 43.5% for at least 12 of the 14 laboratory measurements). Breakdown of country level data is shown in Table 1.

**Table 1.** Sites contributing data to the consortium.

Healthcare System	Acronym	City	Country	Population	Hospitals	Beds	Inpatient discharges/year
Assistance Publique - Hôpitaux de Paris	APHP	Paris	France	Adult & Pediatric	39	20,098	1,375,538
Bordeaux University Hospital	FRBDX	Bordeaux	France	Adult & Pediatric	3	2,676	130,033
Erlangen University Hospital	UKER	Erlangen	Germany	Adult & Pediatric	1	1,400	65000
Medical Center, University of Freiburg	UKFR	Freiburg	Germany	Adult & Pediatric	1	1,660	71,500
University Medicine Mannheim	UMM	Mannheim	Germany	Adult & Pediatric	1	1,352	50,748
ICSM Pavia Hospital	ICSM1	Pavia	Italy	Adult	1	426	8,616
ICSM Lumezzane/Brescia Hospitals	ICSM5	Lumezzane/ Brescia	Italy	Adult	1	149	1,296
ICSM Milano Hospital	ICSM20	Milan	Italy	Adult	1	200	2,432
Policlinico di Milano	POLIMI	Milan	Italy	Adult & Pediatric	1	900	40,000
ASST Papa Giovanni XXIII Bergamo	HPG23	Bergamo	Italy	Adult & Pediatric	1	1,080	45,000
National University Hospital	NUH	Singapore	Singapore	Adult & Pediatric	1	1,556	100,977
Boston Children's Hospital	BCH	Boston, MA	USA	Pediatric	1	404	28,000
Beth Israel Deaconess Medical Center	BIDMC	Boston, MA	USA	Adult	1	673	40,752
Children's Hospital of Philadelphia	CHOP	Philadelphia, PA	USA	Pediatric	1	564	25,406
University of Kansas Medical Center	KUMC	Kansas City, KS	USA	Adult & Pediatric	1	794	54,659
Mayo Clinic	MAYOC	Rochester, MN	USA	Adult & Pediatric	1	2,059	100,000
Mass General Brigham (Partners Healthcare)	MGB	Boston, MA	USA	Adult & Pediatric	10	3,418	163,521
Medical University of South Carolina	MUSC	Charleston, SC	USA	Adult & Pediatric	8	1,600	55,664
University of Pennsylvania	UPenn	Philadelphia, PA	USA	Adult	5	2,469	118,188
University of California, LA	UCLA	Los Angeles, CA	USA	Adult &	2	786	40,526

		CA		Pediatric			
University of Michigan	UMICH	Ann Arbor, MI	USA	Adult & Pediatric	3	1,000	49,008
University of North Carolina at Chapel Hill	UNC	Chapel Hill, NC	USA	Adult & Pediatric	11	3,095	52,000
UT Southwestern Medical Center	UTSW	Dallas, TX	USA	Adult	1	608	26,905
				<b>Total</b>	<b>96</b>	<b>45,352</b>	<b>2,444,792</b>

Demographic breakdown by age and gender is shown in Figure 2. Age distribution was different across countries and consistent with previously identified patterns. In particular, patients from Italy were more commonly over the age of 70 relative to other countries.<sup>23</sup> US institutions, despite representing a large number of active infections, had the lowest percentage of elderly patients diagnosed with COVID-19. Germany, with its three included hospitals and relatively small number of patients, was more similar to the US and had an increased number of male patients in the 50-59 age group.



**Figure 3.** Normalized change (relative to previous day) of new cases reported by 4CE contributors compared to new cases collected by the Johns Hopkins Center for Systems Science and Engineering (JHU CSSE) by country over time.

We were able to capture the total number of identified new cases by site and date. To normalize across sites with a small number of total cases, we generated the rate of growth of total admissions by country. In Figure 3, we compared those values with Johns Hopkins-curated data<sup>24</sup> over time. Rates of growth as extracted from EHR data from our sites were similar to population-level findings during the month of March. Of note, national data from Singapore did not track with EHR numbers. Our one site in Singapore had a small number of patients and is less likely to be representative of all hospitals in the country.

**Despite the limited amount of diagnosis code data submitted, consistent identified symptoms were recorded across institutions.** Most common codes involved respiratory symptoms and infections. Cough, dyspnea, and hyperpyrexia were also commonly identified. Although there were differences in the most

common codes used by sites, symptoms were consistent with previous prospective studies.<sup>25,26</sup> Rates and types of presenting symptoms were similar in the pediatric population.

## Laboratory Value Trajectories

Our initial data extraction comprised 14 laboratory test values that have been strongly associated with poor outcomes in COVID-19 patients in previous publications. The set encompassed markers of cardiac, renal, hepatic, and immune dysfunction. Laboratory trajectories of each hospital at the population level are presented online at <https://covidclinical.net>. Given limitations of data harmonization and space, we focus on 5 laboratory trajectories that encompass immune, hepatic, coagulation, and renal function. Trajectory data were remarkably consistent for most institutions at day 1 (day when biological test positive) with growing differences with continued hospitalization. Extensive data harmonization was performed, but we must emphasize that data from each day represents a potentially different population as patients are discharged, die, or laboratory studies are no longer performed. Data values from each hospital were an average of all studied patients a specified number of days after diagnosis.

**Laboratory values reflected relatively moderate disease severity on presentation.** Initial laboratory values were abnormal for all patients but were not indicative for organ failure. Major abnormal elevations were noted in C-reactive protein (CRP) and D-dimer on the day of diagnosis. As the number of days from diagnosis progressed, laboratory values collected significantly worsened; remaining patients who were not discharged or had died had worsening values. For nearly all 14 tests, trends toward progressively abnormal values were consistent with worsening disease as inpatient stays continued. Most importantly, the initial values and trajectories were highly consistent with previous findings in studies from China.<sup>19,27</sup>

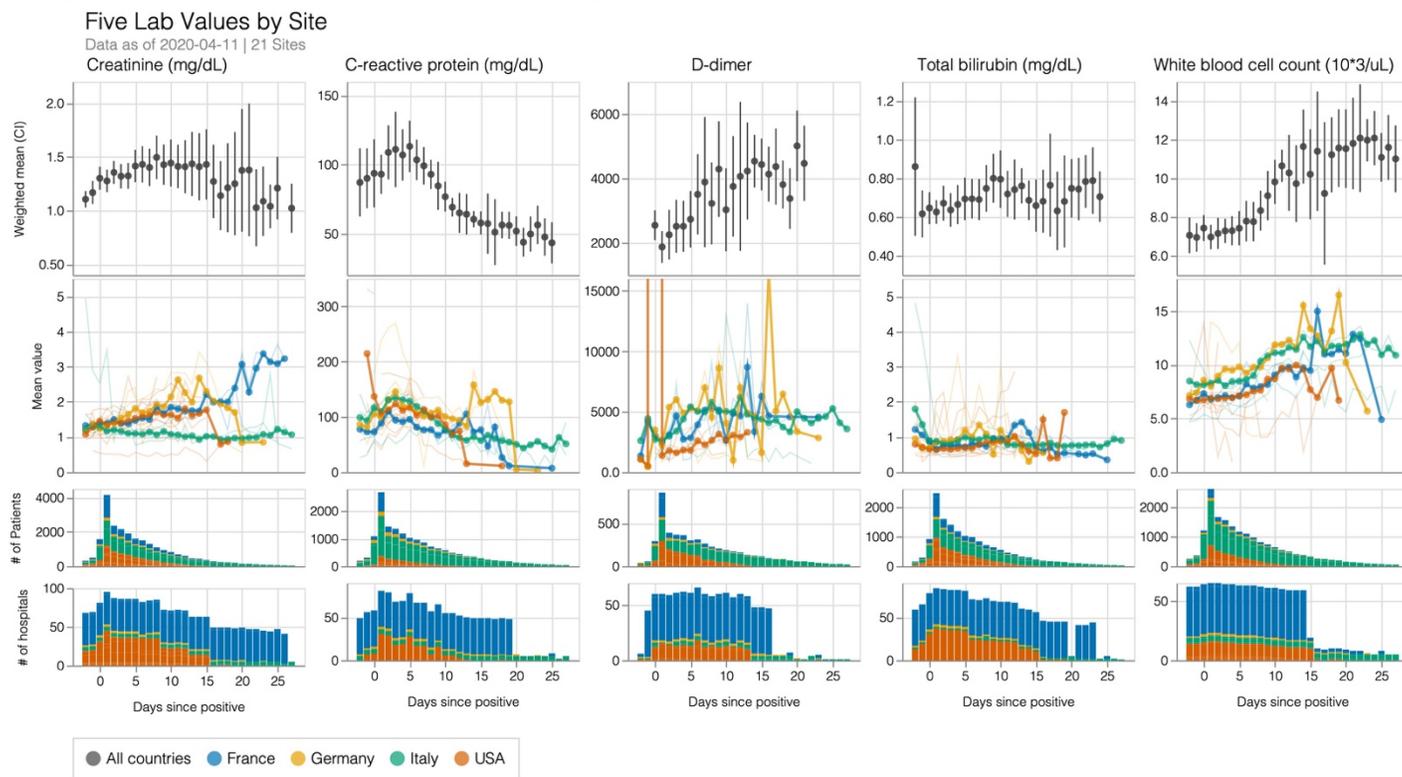
Creatinine, a measure of renal function and the most commonly performed laboratory test in our data set, was divergent over time across sites. This is consistent with an increased proportion of ill patients with significant acute kidney injury over time. Some sites did not have changes in creatinine over time. Hospitals in Italy, in particular, did not see a dramatic rise in creatinine in their hospitalized population. Conversely, the small number of French and German patients remaining in the hospital for two weeks had clear signs of acute kidney injury. This may represent a high mortality near the beginning of the hospitalization at Italian hospitals, severe right censoring of remaining patients, or a difference in practice.

Total bilirubin, a measure of conjugation and function by the liver, was initially normal across most sites and showed increases—consistent with other hepatic laboratory tests—among persistently hospitalized patients. The other hepatic laboratory measurements, ALT and AST, were divergent across institutions and showed a more significant perturbation (see <https://covidclinical.net>). Hepatic impairment was not present in most patients on presentation and total bilirubin was only mildly elevated with continued hospitalization.

On average, white blood cell count (WBC), a measure of immune response, was within normal limits on presentation. Patients who remained in the hospital and survived had increasing WBCs over time without severe leukocytosis.<sup>27</sup> Lymphocyte and neutrophil count trajectories can be seen on the website. Procalcitonin and LDH were not commonly tested in the total patient population, but results can also be seen online.

C-reactive protein (CRP), a measure of systemic inflammation, was notably elevated on presentation for all patients in the cohort with a very narrow confidence interval, consistent with previous findings.<sup>19</sup> Although it is of unclear importance, populations of patients who remained in the hospital, survived, and had ongoing laboratory testing showed improvements over time. Interestingly, despite a decreasing trajectory during the first week, a mild leukocytosis is observed in counterbalance during the second week. The implication may be that CRP is not predictive of ongoing hospitalization or CRP is being checked for patient populations where the laboratory is more commonly improving.

D-dimer, an acute phase reactant and measure of coagulopathy, was elevated across institutions and countries at presentation. It rose consistently in all populations who continued to be hospitalized with the disease. This is consistent with multiple studies that have shown a prothrombotic element to the disease. Most importantly, changes were consistent across all sites and highly abnormal.

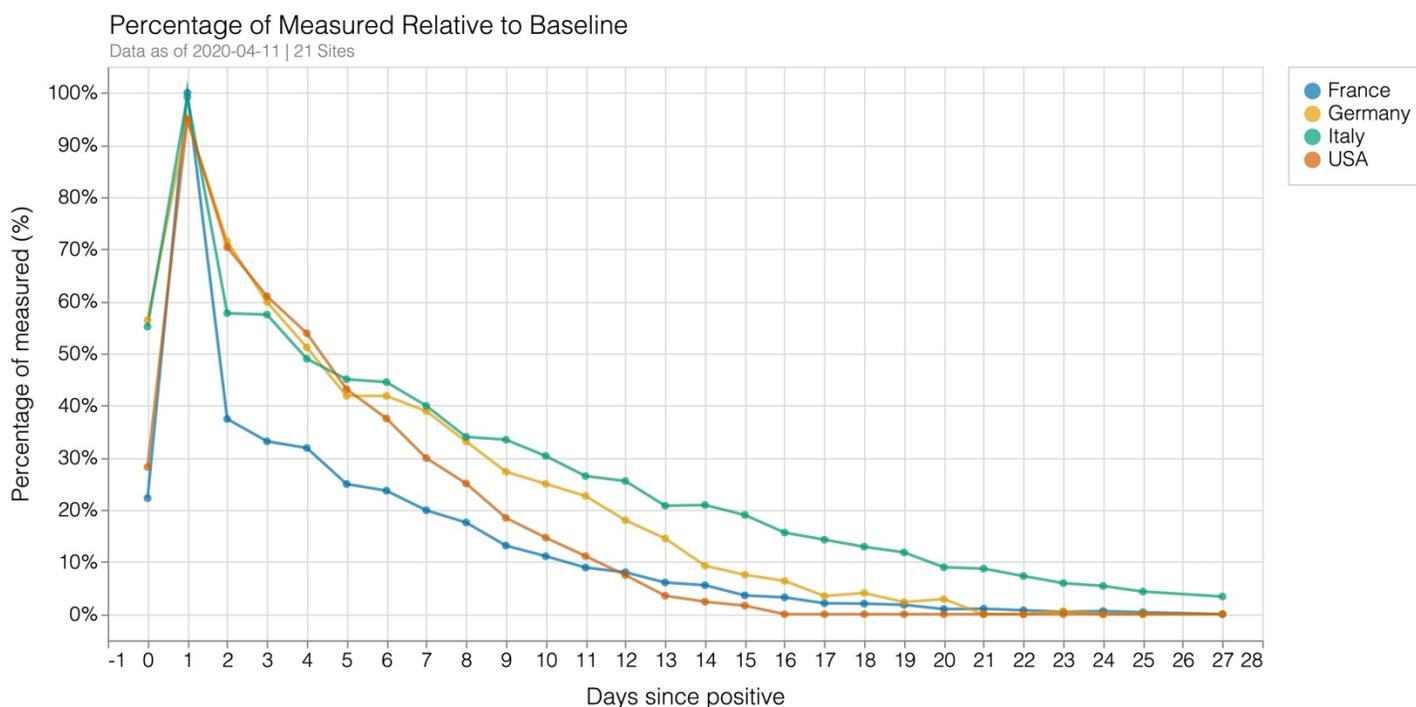


**Figure 4.** Laboratory tests representative of renal function (creatinine), systemic inflammation (C-reactive protein), coagulopathy (D-dimer), liver function (total bilirubin), and immune response (white blood cell count) visualized relative to date of diagnosis of COVID-19. The top row shows weighted means and 95% confidence intervals across all patients. The second row shows unweighted country- (thick lines) and site-level (thin lines) means. The third and fourth rows show the number of patients and sites, respectively, contributing laboratory tests of each type on a given day.

## Data Attrition

There was a large drop in the number of laboratory tests performed after the first day. Drop off in tests performed could be a result of death, length of stay, or lack of interest in further data collection by the clinical team. From the maximum number of laboratory tests consistently checked on the first day after diagnosis, there was a rapid tapering in frequency of laboratory tests checked. These changes were particularly pronounced in

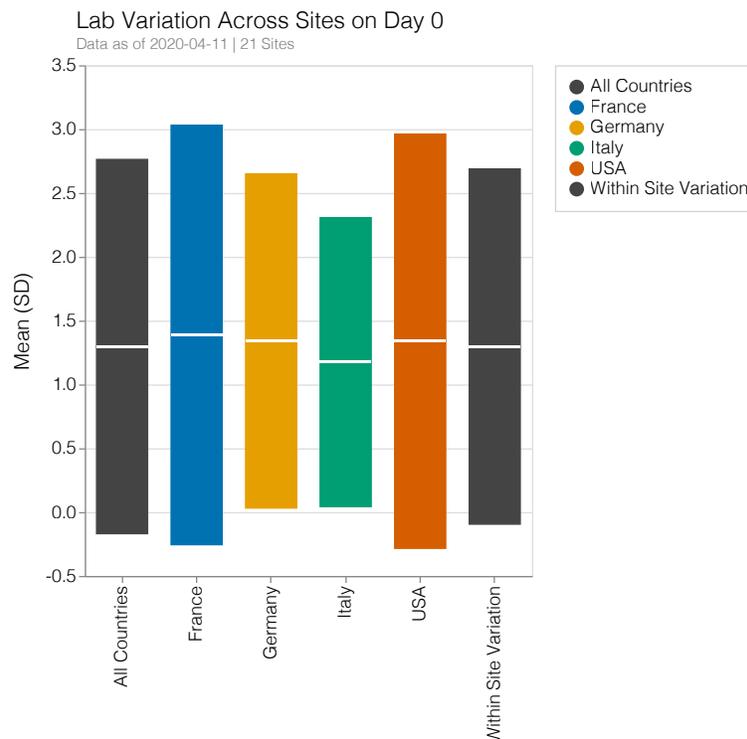
Italy and France. We identified the number of days until the number of tests checked were 20% of their initial maximum value. Values for laboratory study for each day are presented on <https://covidclinical.net>. Results varied for each laboratory value and site. There was no obvious country level pattern. Given that several of these tests, such as creatinine, are commonly checked nearly every day in ill patients, the implication is that patients were censored from the laboratory results because of discharge or death or changing practice pattern. Thus, for the purposes of this paper, we focused on trends in creatinine laboratory testing. We normalized the number of tests performed by day to the total performed on day 1. We then looked at the day when the number of tests performed was 20% of the maximum number performed for each site. **For creatinine, for example, a drop-off in testing occurred between day 7 and 15 across institutions.** Most patients who survived were likely discharged within this time frame or managed with much less monitoring. No country level differences were obvious for this test. Further results can be found online.



**Figure 5.** Drop-off in laboratory tests reported for creatinine relative to the day with the largest number of laboratory tests reported.

## Differences at Admission

There was greater between-hospital variation for laboratory test performance than between-country variation (Figure 6). At the time of diagnosis, there was significant variation between countries and between the hospitals in a specific country. There was no obvious signature presentation for a country for an individual laboratory value. For example, creatinine was a commonly performed laboratory study within a day of diagnosis. The overall standard deviation (SD) for test result values across countries was 1.47 while the SD within sites was 1.39. Standard deviation for countries was 1.64, 1.31, 1.13, and 1.62 within France, Germany, Italy, and the US, respectively. France is a special case as 39 hospitals were reported together by AP-HP and then compared with three hospitals in Bordeaux. This is an important finding that could suggest that laboratory values, as individual results, will not be able to fully explain the mortality differences between countries.



**Figure 6.** Laboratory variation across countries and within sites for laboratory tests performed within a day of diagnosis. Values for mean value and standard deviation (SD) in creatinine level are shown. Large variations exist within sites and are often larger than the between country variation. The overall SD across countries was 1.47 while the SD within sites was 1.39. Standard deviation for countries was 1.64, 1.31, 1.13, and 1.62 within France, Germany, Italy, and the US, respectively.

## Discussion

A rapid mobilization of a multi-national consortium was able to harmonize and integrate data across 5 countries and 3 continents in order to begin to answer questions about comparative care of COVID-19 patients and opportunities for international learning. In just over 2 weeks, the group was able to define a question and data model, perform data extraction and harmonization, evaluate the data, and create a site for public evaluation of site level data. We aggregated EHR data from 96 hospitals, covering a total of 27,927 patients seen in these hospitals for COVID-19. In doing so, we relied upon prior investments made by various governments and institutions in turning the byproducts of clinical documentation into data useful for a variety of operational and scientific tasks, using i2b2 and other implementations, as documented in Table 1. Most importantly, at each site there were bioinformatics experts who understood both the technical characteristics of the data and their clinical relevance.

Using automated data extraction methods, we were able to show results consistent with country-level demographic and epidemiological differences identified in the literature. Rates of total case rise in our study was consistent with international tracking sites.<sup>24</sup> Age breakdown, with Italian sites reporting a larger proportion of older patients, was also reflective of recent publicly available resources.<sup>23</sup>

We were able to show that laboratory trajectories across many hospitals could be collected and were concordant with findings from the literature. In truth, the findings generate more questions than they answer; the ability to see consistencies that spanned many countries indicated that the pathophysiology of this disease is shared across countries, and that demographics and care characteristics will have a significant effect on outcomes. As an example, the fall of CRP among those who continued to be hospitalized with a continued rise in d-dimer could suggest that d-dimer may be more closely related to persistent illness than CRP. The limits of our data collection method, where these results are not tied to the patient level and can be associated across populations, highlights the need for caution with any conclusion related to changes in laboratory levels over time.

Perhaps most importantly, our study did not show a unique laboratory signature at the country level at the time of diagnosis. Researchers around the world have been closely following the rapid spread of COVID-19 and its high mortality rate in certain countries such as Italy. One possible explanation would be that patients who presented to hospitals in Italy did so at a much more advanced stage of disease. Our results do not support this idea. There was as much in-hospital and between-hospital variation as between countries.

The average of laboratory values at presentation did not indicate major organ failure. This may be due to a large proportion of healthier patients than those with advanced disease. Of course, respiratory failure could not be tracked within the limits of our data set.

There were both logistic and data interoperability lessons that were very important to the success of the project and will be critical for future efforts. Logistically, to maximize the timeliness of this consortium's first collaboration around COVID-19, we deliberately aggregated the data to expedite the institutional review board (IRB) process at each institution for such data sharing. This thereby constrained our analyses to count, rather than patient-level, data. While the latter would be optimal for deep analysis and identification of subtle patterns and perturbations of clinical courses, we feel that aggregated count data provide valuable information on the clinical course even as we seek IRB permission for analyses that are at the patient level.

Regarding interoperability, large variations in units and data presentation required extensive data harmonization. The use of LOINC codes allowed for more rapid data extraction, but often institutions did not have internal mappings from their laboratory tests to LOINC codes. Manual interpretation of laboratory value descriptions were sometimes necessary. In future iterations, sites will perform unit conversion and ensure data consistency by presenting reference ranges and example data for a first-pass check of data at the site. Variations in ICD coding as well generation of codes used made code harmonization difficult. Frequencies of presenting codes were useful to show similar patterns to previous literature, but the current set of codes were too sparse for any further meaningful analysis. Future iterations of this project would encompass a much longer data capture timeline and would ensure comprehensive code collection across all sites.

In addition, data alignment by a metric that indicates clinical status is necessary to better establish outcomes. Using day of diagnosis as an alignment strategy does not allow for clear identification of causes for temporal patterns. Similarly, outcomes need to be selected that represent clinically meaningful endpoints secondary to this initial data alignment. One reason for this difficulty was that Identification of level of care was not easily performed. Accordingly, it was not easy to follow patients into and out of ICUs at the site level and ICU data was not reliable.

Our group, the Consortium for Clinical Characterization of COVID-19 by EHR (4CE), is one of hundreds of US efforts (many of which are listed at [HealthIT.gov](https://www.healthit.gov)) that are working hard to aggregate and curate data to inform clinicians, scientists, policy makers and the general public. Additionally, networks of healthcare organizations such as the CTSA's ACT network<sup>28</sup> and PCORnet<sup>29</sup> are working with federal authorities to obtain data-driven population-level insights. Similar initiatives are active in the other countries participating in 4CE, including the German Medical Informatics Initiative.<sup>30</sup> Disease-specific and organ-specific COVID-19 research collectives are also assembling, including ones for cancers (<https://ccc19.org>), inflammatory bowel disease (<https://covidibd.org>), and rheumatology,<sup>31</sup> among many others. The WHO maintains a directory of worldwide research efforts on COVID-19 including clinical data collection.<sup>32</sup> Finally, there are dozens of patient self-reporting apps with hundreds of thousands of users worldwide that provide perspectives on the clinical course of the infection outside hospitals.

There are a multitude of limitations to this study not least of which is that it is observational and subject to a variety of biases with perhaps the most severe being that its data are limited to those patients who were seen at or admitted to hospitals, due to severity of illness or other possibly biasing characteristics. Limitations also include heavy right censoring where patient absence can be due to death or discharge, variations in ICD annotations for conditions existing prior to the COVID-19-related admission, delays in updating billing codes or in uploading EHR data to the local analytic data repository. Furthermore, potentially confounding interactions between comorbidities, chronic diseases and their treatments and lifestyle or exposures were not taken into consideration. Again, because of these limitations we were careful to avoid making more than the basic and descriptive conclusions. Over the coming weeks we will be working on quantifying these biases and adjusting for them, if we can. This will include adding data types as well as disaggregating the data to the patient level if and when permitted by IRBs. For the present, with the current limited knowledge of the clinical course of patients suffering from COVID-19, these results add to this small knowledge-base. Our paper strikingly shows the power of harmonized data extraction from EHRs to rapidly study pandemics like COVID-19.

We invite others to join the 4CE consortium by sending a note to [4CE@i2b2foundation.org](mailto:4CE@i2b2foundation.org).

## References

1. Centers for Disease Control and Prevention. Preliminary Estimates of the Prevalence of Selected Underlying Health Conditions Among Patients with Coronavirus Disease 2019 United States, February 12–March 28, 2020. *MMWR Morb. Mortal. Wkly. Rep.* **69**, 382–386 (2020).
2. Gupta, A. H. Does Covid-19 Hit Women and Men Differently? U.S. Isn't Keeping Track. *The New York Times* (2020).
3. Bonow, R. O., Fonarow, G. C., O'Gara, P. T. & Yancy, C. W. Association of Coronavirus Disease 2019 (COVID-19) With Myocardial Injury and Mortality. *JAMA Cardiol* (2020) doi:10.1001/jamacardio.2020.1105.
4. Ahmadpoor, P. & Rostaing, L. Why the immune system fails to mount an adaptive immune response to a Covid-19 infection. *Transpl. Int.* (2020) doi:10.1111/tri.13611.
5. Thachil, J. The versatile heparin in COVID-19. *J. Thromb. Haemost.* (2020) doi:10.1111/jth.14821.
6. Jin, M. & Tong, Q. Rhabdomyolysis as Potential Late Complication Associated with COVID-19. *Emerg. Infect. Dis.* **26**, (2020).

7. Zhang, C., Shi, L. & Wang, F. S. Liver injury in COVID-19: management and challenges. *Lancet Gastroenterol Hepatol* (2020) doi:10.1016/S2468-1253(20)30057-1.
8. Pan, X. W. *et al.* Identification of a potential mechanism of acute kidney injury during the COVID-19 outbreak: a study based on single-cell transcriptome analysis. *Intensive Care Med.* (2020) doi:10.1007/s00134-020-06026-1.
9. Lippi, G. & Plebani, M. Procalcitonin in patients with severe coronavirus disease 2019 (COVID-19): A meta-analysis. *Clin. Chim. Acta* **505**, 190–191 (2020).
10. Murphy, S. N. *et al.* Serving the enterprise and beyond with informatics for integrating biology and the bedside (i2b2). *J. Am. Med. Inform. Assoc.* **17**, 124–130 (2010).
11. Gutierrez-Sacristan, A. *et al.* Rcupcake: an R package for querying and analyzing biomedical data through the BD2K PIC-SURE RESTful API. *Bioinformatics* **34**, 1431–1432 (2018).
12. Mandl, K. D. *et al.* The Genomics Research and Innovation Network: creating an interoperable, federated, genomics learning system. *Genet. Med.* **22**, 371–380 (2020).
13. Mandl, K. D. *et al.* Scalable Collaborative Infrastructure for a Learning Healthcare System (SCILHS): architecture. *J. Am. Med. Inform. Assoc.* **21**, 615–620 (2014).
14. McMurry, A. J. *et al.* SHRINE: enabling nationally scalable multi-site disease studies. *PLoS One* **8**, e55811 (2013).
15. Weber, G. M. *et al.* The Shared Health Research Information Network (SHRINE): a prototype federated query tool for clinical data repositories. *J. Am. Med. Inform. Assoc.* **16**, 624–630 (2009).
16. Visweswaran, S. *et al.* Accrual to Clinical Trials (ACT): A Clinical and Translational Science Award Consortium Network. *JAMIA Open* **1**, 147–152 (2018).
17. Kohane, I. S., Churchill, S. E. & Murphy, S. N. A translational engine at the national scale: informatics for integrating biology and the bedside. *J. Am. Med. Inform. Assoc.* **19**, 181–185 (2012).
18. Bodenreider, O., Cornet, R. & Vreeman, D. J. Recent Developments in Clinical Terminologies - SNOMED CT, LOINC, and RxNorm. *Yearb. Med. Inform.* **27**, 129–139 (2018).
19. Lippi, G. & Plebani, M. Laboratory abnormalities in patients with COVID-2019 infection. *Clinical Chemistry and Laboratory Medicine (CCLM)* vol. 0 (2020).
20. i2b2: Informatics for Integrating Biology & the Bedside. <https://www.i2b2.org/>.
21. CTSA ACT Consortium. *CTSA ACT Network i2b2 and SHRINE Ontology with 1-1 SHRINE Adapter Mapping file.* (Github).
22. Bodenreider, O. The Unified Medical Language System (UMLS): integrating biomedical terminology. *Nucleic Acids Res.* **32**, D267-70 (2004).
23. Onder, G., Rezza, G. & Brusaferro, S. Case-Fatality Rate and Characteristics of Patients Dying in Relation to COVID-19 in Italy. *JAMA* (2020) doi:10.1001/jama.2020.4683.
24. Dong, E., Du, H. & Gardner, L. An interactive web-based dashboard to track COVID-19 in real time. *Lancet Infect. Dis.* (2020) doi:10.1016/S1473-3099(20)30120-1.
25. Huang, C. *et al.* Clinical features of patients infected with 2019 novel coronavirus in Wuhan, China. *Lancet* **395**, 497–506 (2020).
26. Guan, W.-J. *et al.* Clinical Characteristics of Coronavirus Disease 2019 in China. *N. Engl. J. Med.* (2020) doi:10.1056/NEJMoa2002032.
27. Wynants, L. *et al.* Prediction models for diagnosis and prognosis of covid-19 infection: systematic review and critical appraisal. *BMJ* **369**, m1328 (2020).
28. Shah, M. R. *et al.* Early Vision for the CTSA Program Trial Innovation Network: A Perspective from the National Center for Advancing Translational Sciences. *Clin. Transl. Sci.* **10**, 311–313 (2017).
29. Fleurence, R. L. *et al.* Launching PCORnet, a national patient-centered clinical research network. *J. Am. Med. Inform. Assoc.* **21**, 578–582 (2014).
30. Semler, S. C., Wissing, F. & Heyder, R. German Medical Informatics Initiative. *Methods Inf. Med.* **57**, e50–e56 (2018).
31. Robinson, P. C. & Yazdany, J. The COVID-19 Global Rheumatology Alliance: collecting data in a pandemic. *Nat. Rev. Rheumatol.* (2020) doi:10.1038/s41584-020-0418-0.

32. Global research on coronavirus disease (COVID-19). <https://www.who.int/emergencies/diseases/novel-coronavirus-2019/global-research-on-novel-coronavirus-2019-ncov>.

**Affiliations:**

First Name	Initial	Last Name	Primary Affiliation
Katie		Kirchoff	Biomedical Informatics Center, Medical University of South Carolina
Jean		Craig	Biomedical Informatics Center, Medical University of South Carolina
Jihad		Obeid	Biomedical Informatics Center, Medical University of South Carolina
Mauro		Bucalo	BIOMERIS (BIOMedical Research Informatics Solutions)
Romain		Griffier	Bordeaux University Hospital
Bertrand		Moal	Bordeaux University Hospital
Vianney		Jouhet	Bordeaux University Hospital / ERIAS - Inserm U1219 BPH
Sébastien		Cossin	Bordeaux University Hospital / ERIAS - Inserm U1219 BPH
Detlef		Kraska	Center for Medical Information and Communication Technology, University Hospital Erlangen
Piotr		Sliz	CHIP, Boston Children's Hospital
Damien		Leprovost	Clevy.io
Judith		Leblanc	Clinical Research Unit, Saint Antoine Hospital, APHP Greater Paris University Hospital
Mohamad		Daniar	Clinical Research Informatics, Boston Children's Hospital
Patricia		Martel	Clinical Research Unit, Paris Saclay, APHP Greater Paris University Hospital
Kenneth	D	Mandl	Computational Health Informatics Program, Boston Children's Hospital
Batsal		Devkota	Department of Biomedical and Health Informatics, Children's Hospital of Philadelphia
Isaac	S	Kohane	Department of Biomedical Informatics, Harvard Medical School
Paul		Avillach	Department of Biomedical Informatics, Harvard Medical School
Griffin	M	Weber	Department of Biomedical Informatics, Harvard Medical School
Nathan	P	Palmer	Department of Biomedical Informatics, Harvard Medical School
Nils		Gehlenborg	Department of Biomedical Informatics, Harvard Medical School
Sehi		L'Yi	Department of Biomedical Informatics, Harvard Medical School
Mark	S	Keller	Department of Biomedical Informatics, Harvard Medical School
Arnaud		Serret-Larmande	Department of Biomedical Informatics, Harvard Medical School
Alba		Gutiérrez-Sacristán	Department of Biomedical Informatics, Harvard Medical School
Amelia	LM	Tan	Department of Biomedical Informatics, Harvard Medical School
Brett	K	Beaulieu-Jones	Department of Biomedical Informatics, Harvard Medical School
Anne Sophie		Jannot	Department of Biomedical Informatics, HEGP, APHP Greater Paris University Hospital
Anita		Burgun	Department of Biomedical Informatics, HEGP, APHP Greater Paris University Hospital
John	H	Holmes	Department of Biostatistics, Epidemiology, and Informatics, University of Pennsylvania Perelman School of Medicine
Danielle	L	Mowery	Department of Biostatistics, Epidemiology, and Informatics/Institute for Biomedical Informatics, University of Pennsylvania Perelman School of Medicine
Riccardo		Bellazzi	Department of Electrical Computer and Biomedical Engineering, University of Pavia, Italy and IRCCS ICS Maugeri, Italy
Lemuel	R	Waitman	Department of Internal Medicine, Division of Medical Informatics
Lav	P	Patel	Department of Internal Medicine, Division of Medical Informatics, University Of Kansas Medical Center
Hans	U	Prokosch	Department of Medical Informatics, University of Erlangen-Nürnberg
Douglas	S	Bell	Department of Medicine, David Geffen School of Medicine at UCLA

Robert	W	Follett	Department of Medicine, David Geffen School of Medicine at UCLA
Douglas	A	Murad	Department of Medicine, David Geffen School of Medicine at UCLA
Jeffrey	G	Klann	Department of Medicine, Massachusetts General Hospital
Shawn	N	Murphy	Department of Neurology, Massachusetts General Hospital
Alberto		Zambelli	Department of Oncology, ASST Papa Giovanni XXIII, Bergamo
David	A	Hanauer	Department of Pediatrics, University of Michigan Medical School
Gilbert	S	Omenn	Dept of Computational Medicine & Bioinformatics, U of Michigan
Arthur		Mensch	ENS, PSL University
Thomas		Ganslandt	Heinrich-Lanz-Center for Digital Health, University Medicine Mannheim, Heidelberg University
Tobias		Gradinger	Heinrich-Lanz-Center for Digital Health, University Medicine Mannheim, Heidelberg University
Julien		Champ	INRIA Sophia-Antipolis – ZENITH team, LIRMM, Montpellier, France
Jason	H	Moore	Institute for Biomedical Informatics, University of Pennsylvania
Christian		Haverkamp	Institute of Digitalization in Medicine, Faculty of Medicine and Medical Center, University of Freiburg, Germany
Martin		Boeker	Institute of Medical Biometry and Statistics, Medical Center, University of Freiburg
Valentina		Tibollo	IRCCS ICS Maugeri, Pavia
Alberto		Malovini	IRCCS ICS Maugeri, Pavia
Luca		Chiovato	IRCCS ICS Maugeri, Pavia and Department of Internal Medicine and medical Therapy, University of Pavia
Kee Yuan		Ngiam	National University Health Systems Singapore
Robert	L	Bradford	North Carolina Translational and Clinical Sciences (NC TraCS) Institute, UNC Chapel Hill
Emily		Schiver	Penn Medicine, Data Analytics Center
James		Cimino	UAB Informatics Institute
Luigia		Scudeller	Scientific Direction, IRCCS Ca' Granda Ospedale Maggiore Policlinico di Milano
Loic		Esteve	SED/SIERRA, Inria Centre de Paris
Jill Jen		Vie	SequeL, Inria Lille
Mélodie		Bernaux	Strategy and transformation department, APHP Greater Paris University Hospital
Alexandre		Gramfort	Université Paris-Saclay, Inria, CEA
Olivier		Grisel	Université Paris-Saclay, Inria, CEA
Guillaume		Lemaitre	Université Paris-Saclay, Inria, CEA
Thomas		Moreau	Université Paris-Saclay, Inria, CEA
Demian		Wassermann	Université Paris-Saclay, Inria, CEA
Charlotte		Caucheteux	Université Paris-Saclay, Inria, CEA
Gael		Varoquaux	Université Paris-Saclay, Inria, CEA, Montréal Neurological Institute, McGill University
Antonio		Bellasi	UOC Ricerca, Innovazione e Brand reputation, ASST Papa Giovanni XXIII, Bergamo
Vincent		Benoit	WIND Department APHP Greater Paris University Hospital
Romain		Bey	WIND Department APHP Greater Paris University Hospital
Nicolas		Paris	WIND Department APHP Greater Paris University Hospital
Patricia		Serre	WIND Department APHP Greater Paris University Hospital
Nina		Orlova	WIND Department APHP Greater Paris University Hospital
Julien		Dubiel	WIND Department APHP Greater Paris University Hospital
Martin		Hilka	WIND Department APHP Greater Paris University Hospital

Stéphane		Bréant	WIND Department APHP Greater Paris University Hospital
Arnaud		Sandrin	WIND Department APHP Greater Paris University Hospital
Elisa		Salamanca	WIND Department APHP Greater Paris University Hospital
Sylvie		Cormont	WIND Department APHP Greater Paris University Hospital
Christel		Daniel	WIND Department APHP Greater Paris University Hospital UMRS1142 INSERM
Nicolas		Griffon	WIND Department APHP Greater Paris University Hospital UMRS1142 INSERM
Nicolas		Griffon	WIND Department APHP Greater Paris University Hospital UMRS1142 INSERM