



HAL
open science

Potential pitfalls in the use of real-world data for studying long COVID

Harrison G Zhang, Jacqueline P Honerlaw, Monika Maripuri, Malarkodi Jebathilagam Samayamuthu, Brendin R Beaulieu-Jones, Huma S Baig, Sehi L'yi, Yuk-Lam Ho, Michele Morris, Vidul Ayakulangara Panickan, et al.

► **To cite this version:**

Harrison G Zhang, Jacqueline P Honerlaw, Monika Maripuri, Malarkodi Jebathilagam Samayamuthu, Brendin R Beaulieu-Jones, et al. Potential pitfalls in the use of real-world data for studying long COVID. *Nature Medicine*, 2023, 29 (5), pp.1040-1043. 10.1038/s41591-023-02274-y . hal-04575574

HAL Id: hal-04575574

<https://hal.science/hal-04575574>

Submitted on 6 Jun 2024

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Published in final edited form as:

Nat Med. 2023 May ; 29(5): 1040–1043. doi:10.1038/s41591-023-02274-y.

†Corresponding author: Gabriel A Brat, MD, MPH. Department of Biomedical Informatics, Harvard Medical School, Boston, MA, United States.

*These authors contributed equally

The Consortium for Clinical Characterization of COVID-19 by EHR (4CE)

James R Aaron MHA³⁹, Giuseppe Agapito PhD⁴⁰, Adem Albayrak⁴¹, Giuseppe Albi MS²⁸, Mario Alessiani MD, FACS⁴², Anna Alloni PhD³⁷, Danilo F Amendola MSc⁴³, François Angoulyant MD, PhD⁴⁴, Li L.L.J Anthony⁴⁵, Bruce J Aronow PhD³², Fatima Ashraf MS⁴⁶, Andrew Atz MD⁴⁷, Paul Avillach MD, PhD¹, Paula S Azevedo MD, PhD⁴⁸, James Balshi⁴⁹, Brett K Beaulieu-Jones PhD¹, Douglas S Bell²³, Antonio Bellasi MD, PhD⁵⁰, Riccardo Bellazzi MS, PhD²⁸, Vincent Benoit PhD²⁴, Michele Beraghi MS⁵¹, José Luis Bernal-Sobrino MS⁵², Mélodie Bernaux⁵³, Romain Bey²⁴, Surbhi Bhatnagar PhD³², Alvar Blanco-Martínez MS⁵², Clara-Lea Bonzel MSc¹, John Booth MSc⁵⁴, Silvano Bosari Prof.³⁵, Florence T Bourgeois MD, MPH¹¹, Robert L Bradford⁵⁵, Gabriel A Brat MD¹, Stéphane Bréant⁵⁶, Nicholas W Brown MEng¹, Raffaele Bruno MD⁵⁷, William A Bryant PhD⁵⁴, Mauro Bucalo MS³⁷, Emily Bucholz MD, PhD, MPH⁵⁸, Anita Burgun⁵⁹, Tianxi Cai ScD¹, Mario Cannataro M.Sc.⁶⁰, Aldo Carmona⁶¹, Charlotte Caucheteux⁶², Julien Champ⁶³, Jin Chen PhD⁶⁴, Krista Y Chen BS⁶⁵, Luca Chiovato MD, PhD³¹, Lorenzo Chiudinelli PhD⁶⁶, Kelly Cho PhD, MPH²⁹, James J Cimino MD⁶⁷, Tiago K Colicchio PhD, MBA⁶⁷, Sylvie Cormont⁵⁶, Sébastien Cossin³⁰, Jean B Craig PhD⁶⁸, Juan Luis Cruz-Bermúdez PhD⁵², Jaime Cruz-Rojo MD⁵², Arianna Dagliati MS, PhD², Mohamad Danian MSIS⁶⁹, Christel Daniel⁷⁰, Priyam Das PhD¹, Batsal Devkota⁷¹, Audrey Dionne MD⁵⁸, Rui Duan PhD³, Julien Dubiel⁵⁶, Scott L DuVall PhD⁷², Loïc Esteve⁷³, Hossein Estiri PhD¹⁵, Shirley Fan⁷⁴, Robert W Follett BS²³, Thomas Ganslandt MD⁷⁵, Noelia García-Barrio MS⁵², Lana X Garmire PhD⁷⁶, Nils Gehlenborg¹, Emily J Getzen MS⁷⁷, Alon Geva MD, MPH⁷⁸, Tobias Gradinger MD, BSc⁷⁵, Alexandre Gramfort⁶², Romain Griffier³⁰, Nicolas Griffon⁷⁰, Olivier Grisel⁶², Alba Gutiérrez-Sacristán PhD¹, Larry Han PhD³, David A Hanauer MD, MS⁸, Christian Haverkamp MD⁷⁹, Derek Y Hazard MSc⁸⁰, Bing He PhD⁷⁶, Darren W Henderson BS³⁹, Martin Hilka⁵⁶, Yuk-Lam Ho MPH²², John H Holmes MS, PhD^{9,10}, Chuan Hong PhD^{1,8†}, Kenneth M Huling HS¹, Meghan R Hutch BS⁸², Richard W Issitt DClintP⁵⁴, Anne Sophie Jannot⁸³, Vianney Jouhet MD, PhD³⁰, Ramakanth Kavuluru PhD²⁵, Mark S Keller¹, Chris J Kennedy PhD⁸⁴, Daniel A Key BEng⁵⁴, Katie Kirchoff MSHI⁸⁵, Jeffrey G Klann MEng, PhD¹⁵, Isaac S Kohane MD, PhD¹, Ian D Krantz⁸⁶, Detlef Kraska Dr.⁸⁷, Ashok K Krishnamurthy PhD⁸⁸, Sehi L Yi PhD¹, Trang T Le PhD⁹, Judith Leblanc⁸⁹, Guillaume Lemaitre⁶², Leslie Lenert MD, MS⁶⁸, Damien Leprovost⁹⁰, Molei Liu PhD⁹¹, Ne Hooi Will Loh MBBS⁹², Qi Long PhD⁹³, Sara Lozano-Zahonero PhD²⁰, Yuan Luo PhD⁸², Kristine E Lynch PhD⁷², Sadiqa Mahmood⁴¹, Sarah E Maidlow AA¹², Adeline Makoudjou MD²⁰, Alberto Malovini PhD²⁷, Kenneth D Mandl MD, MPH⁶⁵, Chengsheng Mao PhD⁸², Anupama Maram MS⁹⁴, Patricia Martel⁹⁵, Marcelo R Martins MSc⁹⁶, Jayson S Marwaha MD⁹⁷, Aaron J Masino PhD⁹⁸, Maria Mazzitelli PhD⁹⁹, Arthur Mensch¹⁰⁰, Marianna Milano PhD¹⁰¹, Marcos F Minicucci MD, PhD¹⁰², Bertrand Moal MD, PhD¹³, Taha Mohseni Ahooyi PhD¹⁰³, Jason H Moore PhD¹⁰⁴, Cinta Moraleta MD, PhD¹⁰⁵, Jeffrey S Morris¹⁰⁶, Michele Morris BA⁶, Karyn L Moshal¹⁰⁷, Sajad Mousavi PhD¹, Danielle L Mowery PhD⁹, Douglas A Murad²³, Shawn N Murphy MD, PhD¹⁴, Thomas P Noughton BA¹⁰⁸, Carlos Tadeu Breda Neto⁴³, Antoine Neuraz MD, PhD¹⁶, Jane Newburger MD, MPH⁵⁸, Kee Yuan Ngiam MBBS, FRCS¹⁰⁹, Wanjiku FM Njoroge MD¹¹⁰, James B Norman¹, Jihad Obeid MD, FAMIA⁶⁸, Marina P Okoshi PhD¹⁰², Karen L Olson PhD¹¹¹, Gilbert S. Omenn MD, PhD¹⁹, Nina Orlova⁵⁶, Brian D Ostasiewski BS¹¹², Nathan P Palmer PhD¹, Nicolas Paris⁵⁶, Lav P Patel MS⁷, Miguel Pedrera-Jiménez MS⁵², Emily R Pfaff PhD¹¹³, Ashley C Pfaff MD¹¹⁴, Danielle Pillion MS¹, Sara Pizzimenti MS³⁵, Hans U Prokosch¹¹⁵, Robson A Prudente PhD¹¹⁶, Andrea Prunotto PhD²⁰, Víctor Quirós-González MS⁵², Rachel B Ramoni¹¹⁷, Maryna Raskin⁴¹, Siegfert Rieg MD¹¹⁸, Gustavo Roig-Domínguez MS⁵², Pablo Rojo MD, PhD¹¹⁹, Paula Rubio-Mayo MS⁵², Paolo Sacchi MD⁵⁷, Carlos Sáez PhD¹²⁰, Elisa Salamanca⁵⁶, Malarkodi Jebathilagam Samayamuthu MD⁶, L. Nelson Sanchez-Pinto MD, MBI¹²¹, Arnaud Sandrin⁵⁶, Nandhini Santhanam MSc⁷⁵, Janaina C.C Santos MS¹²², Fernando J Sanz Vidoreta²³, Maria Savino MS¹²³, Emily R Schriver MS¹²⁴, Petra Schubert MPH²², Juergen Schuettler¹²⁵, Luigia Scudeller MD, MSc³⁵, Neil J Sebire MD, FRCPATH⁵⁴, Pablo Serrano-Balazote MD, MS⁵², Patricia Serre⁵⁶, Arnaud Serret-Larmande MD¹²⁶, Mohsin Shah MSc⁵⁴, Zahra Shakeri Hossein Abad PhD¹, Domenick Silvio¹²⁷, Piotr Sliz⁶⁵, Jiyeon Son MD¹²⁸, Charles Sunday¹²⁹, Andrew M South MD, MS³⁸, Anastasia Spiridou PhD⁵⁴, Zachary H. Strasser MD¹⁵, Amelia LM Tan BSc, PhD¹, Bryce W.Q. Tan MBBS⁵, Byorn W.L. Tan MBBS⁵, Suzana E Tanni PhD¹⁰², Deanne M Taylor PhD¹³⁰, Ana I Terriza-Torres MS⁵², Valentina Tibollo MS²⁷, Patric Tippmann MSc⁸⁰, Emma MS Toh³³, Carlo Torti PhD⁹⁹, Enrico M Treccarichi PhD⁹⁹, Yi-Ju Tseng PhD¹³¹, Andrew K Vallejos¹³², Gael Varoquaux¹³³, Margaret E Vella BS¹, Guillaume Verdy MSc¹³, Jill-Jénn Vie¹³⁴, Shyam Visweswaran MD, PhD⁶, Michele Vitacca MD, PhD¹³⁵, Kavishwar B Wagholikar MBBS, PhD³⁶, Lemuel R Waitman¹³⁶, Xuan Wang PhD¹, Demian Wassermann⁶², Griffin M Weber MD, PhD¹, Martin Wolkewitz PhD⁸⁰, Scott Wong⁵, Zongqi Xia MD, PhD⁴, Xin Xiong MS³, Ye Ye BMED, MSPH, PhD⁶, Nadir Yehya MD, MSCE¹³⁷, William Yuan PhD¹, Alberto Zambelli¹³⁸, Harrison G Zhang BA¹, Daniela Zöller PhD²⁰, Valentina Zuccaro MD⁵⁷, Chiara Zucco PhD¹⁰¹

¹Department of Biomedical Informatics, Harvard Medical School, Boston, United States

²Department of Electrical Computer and Biomedical Engineering, University of Pavia, Pavia, Italy.

³Department of Biostatistics, Harvard T.H. Chan School of Public Health, Boston, United States.

⁴Department of Neurology, University of Pittsburgh, Pittsburgh, United States.

⁵Department of Medicine, National University Hospital, Singapore, Singapore, Singapore.

⁶Department of Biomedical Informatics, University of Pittsburgh, Pittsburgh, United States.

⁷Department of Internal Medicine, Division of Medical Informatics, University Of Kansas Medical Center, Kansas City, United States.

⁸Department of Learning Health Sciences, University of Michigan Medical School, Ann Arbor, United States.

⁹Department of Biostatistics, Epidemiology, and Informatics, University of Pennsylvania Perelman School of Medicine, Philadelphia, United States.

¹⁰Institute for Biomedical Informatics, University of Pennsylvania Perelman School of Medicine, Philadelphia, United States.

¹¹Department of Pediatrics, Harvard Medical School, Boston, United States.

¹²Michigan Institute for Clinical and Health Research (MICH) Informatics, University of Michigan, Ann Arbor, United States.

¹³IAM unit, Bordeaux University Hospital, Bordeaux, France.

¹⁴Department of Neurology, Massachusetts General Hospital, Boston, United States.

- 15 Department of Medicine, Massachusetts General Hospital, Boston, United States.
of Paris, Paris, France.
- 17 Department of Biomedical informatics, WiSDM, National University Health Systems Singapore, Singapore, Singapore.
- 18 Department of Anaesthesia, National University Health Systems Singapore, Singapore, Singapore.
- 19 Dept of Computational Medicine & Bioinformatics, Internal Medicine, Human Genetics, and School of Public Health, University of Michigan, Ann Arbor, United States.
- 20 Institute of Medical Biometry and Statistics, Faculty of Medicine and Medical Center, University of Freiburg, Freiburg, Germany.
- 21 Department of Ophthalmology, Mayo Clinic, Rochester, United States.
- 22 Massachusetts Veterans Epidemiology Research and Information Center (MAVERIC), VA Boston Healthcare System, Boston, United States.
- 23 Department of Medicine, David Geffen School of Medicine at UCLA, Los Angeles, United States.
- 24 IT Department, Innovation & Data, APHP Greater Paris University Hospital, Paris, France.
- 25 Division of Biomedical Informatics (Department of Internal Medicine), University of Kentucky, Lexington, United States.
- 26 Department of Preventive Medicine, Northwestern University, Chicago, USA.
- 27 Laboratory of Informatics and Systems Engineering for Clinical Research, Istituti Clinici Scientifici Maugeri SpA SB IRCCS, Pavia, Italy.
- 28 Department of Electrical, Computer and Biomedical Engineering, University of Pavia, Pavia, Italy.
- 29 Population Health and Data Science, MAVERIC, VA Boston Healthcare System, Boston, United States.
- 30 IAM unit, INSERM Bordeaux Population Health ERIAS TEAM, Bordeaux University Hospital / ERIAS - Inserm U1219 BPH, Bordeaux, France.
- 31 Unit of Internal Medicine and Endocrinology, Istituti Clinici Scientifici Maugeri SpA SB IRCCS, Pavia, Italy.
- 32 Departments of Biomedical Informatics, Pediatrics, Cincinnati Children's Hospital Medical Center, University of Cincinnati, Cincinnati, United States.
- 33 Yong Loo Lin School of Medicine, National University of Singapore, Singapore, Singapore.
- 34 Department of Medicine, National University Health Systems Singapore, Singapore, Singapore.
- 35 Scientific Direction, IRCCS Ca' Granda Ospedale Maggiore Policlinico di Milano, Milan, Italy.
- 36 Department of Medicine, Massachusetts General Hospital, Boston, USA.
- 37 BIOMERIS (BIOMedical Research Informatics Solutions), Pavia, Italy.
- 38 Department of Pediatrics-Section of Nephrology, Brenner Children's, Wake Forest School of Medicine, Winston Salem, United States.
- 39 Department of Biomedical Informatics, University of Kentucky, Lexington, United States.
- 40 Department of Legal, Economic and Social Sciences, University Magna Graecia of Catanzaro, Catanzaro, Italy.
- 41 Health Catalyst, INC., Cambridge, United States.
- 42 Department of Surgery, ASST Pavia, Lombardia Region Health System, Pavia, Italy.
- 43 Clinical Research Unit of Botucatu Medical School, São Paulo State University, Clinical Research Unit of Botucatu Medical School, São Paulo State University, Botucatu, Brazil.
- 44 Pediatric emergency Department, Hôpital Necker-Enfants Malades, Assistance Public-Hôpitaux de Paris, Paris, France. ⁴⁵National Center for Infectious Diseases, Tan Tock Seng Hospital, Singapore, Singapore.
- 46 BIG-ARC, The University of Texas Health Science Center at Houston, School of Biomedical Informatics, Houston, United States.
- 47 Department of Pediatrics, Medical University of South Carolina, Charleston, United States.
- 48 Internal Medicine Department, Botucatu Medical School, São Paulo State University, Botucatu, Brazil.
- 49 Department of Surgery, St. Luke's University Health Network, Bethlehem, United States.
- 50 Department of Medicine, Division of Nephrology, Ente Ospedaliero Cantonale, Lugano, Switzerland.
- 51 IT Department, ASST Pavia, Voghera, Italy.
- 52 Health Informatics, Hospital Universitario 12 de Octubre, Madrid, Spain.
- 53 Strategy and Transformation Department, APHP Greater Paris University Hospital, Paris, France.
- 54 Digital Research, Informatics and Virtual Environments (DRIVE), Great Ormond Street Hospital for Children, UK, London, United Kingdom.
- 55 North Carolina Translational and Clinical Sciences (NC TraCS) Institute, UNC Chapel Hill, Chapel Hill, United States.
- 56 IT department, Innovation & Data, APHP Greater Paris University Hospital, Paris, France.
- 57 Division of Infectious Diseases I, Fondazione I.R.C.C.S. Policlinico San Matteo, Pavia, Italy.
- 58 Department of Cardiology, Boston Children's Hospital, Harvard Medical School, Boston, United States.
- 59 Department of Biomedical Informatics, HEGP, APHP Greater Paris University Hospital, Paris, France.
- 60 Department of Medical and Surgical Sciences, Data Analytics Research Center, University Magna Graecia of Catanzaro, Catanzaro, Italy.
- 61 Department of Anesthesia, St. Luke's University Health Network, Bethlehem, United States.
- 62 Université Paris-Saclay, Inria, CEA, Palaiseau, France.
- 63 INRIA Sophia-Antipolis – ZENITH team, LIRMM, Montpellier, France.
- 64 Department of Internal Medicine, University of Kentucky, Lexington, United States.
- 65 Computational Health Informatics Program, Boston Children's Hospital, Boston, United States.
- 66 UOC Ricerca, Innovazione e Brand reputation, ASST Papa Giovanni XXIII, Bergamo, Bergamo, Italy.
- 67 Informatics Institute, University of Alabama at Birmingham, Birmingham, United States.
- 68 Biomedical Informatics Center, Medical University of South Carolina, Charleston, United States.
- 69 Clinical Research Informatics, Boston Children's Hospital, Boston, United States.
- 70 IT department, Innovation & Data (APHP), UMRS1142 (INSERM), APHP Greater Paris University Hospital, INSERM, Paris, France.
- 71 Department of Biomedical and Health Informatics, Children's Hospital of Philadelphia, Philadelphia, United States.

COVID

- ⁷²VA Informatics and Computing Infrastructure, VA Salt Lake City Health Care System, Salt Lake City, United States.
- ⁷³SED/SIERRA, Inria Centre de Paris, Paris, France.
- ⁷⁴Health Information Technology & Services, University of Michigan, Ann Arbor, United States.
- ⁷⁵Heinrich-Lanz-Center for Digital Health, University Medicine Mannheim, Heidelberg University, Mannheim, Germany.
- ⁷⁶Department of Computational Biology and Bioinformatics, University of Michigan, Ann Arbor, United States.
- ⁷⁷Biostatistics, Perelman School of Medicine at the University of Pennsylvania, Philadelphia, United States.
- ⁷⁸Department of Anesthesiology, Critical Care, and Pain Medicine and Computational Health Informatics Program, Boston Children's Hospital, Boston, United States.
- ⁷⁹Institute of Digitalization in Medicine, Faculty of Medicine and Medical Center, University of Freiburg, Freiburg, Germany.
- ⁸⁰Institute of Medical Biometry and Statistics, Institute of Medical Biometry and Statistics, Medical Center, University of Freiburg, Freiburg, Germany.
- ⁸¹Department of Biostatistics and Bioinformatics, Duke University, Durham, United States.
- ⁸²Department of Preventive Medicine, Northwestern University, Chicago, United States.
- ⁸³Department of Biomedical Informatics, HEGP, APHP Greater Paris University Hospital, Paris, France.
- ⁸⁴Center for Precision Psychiatry, Massachusetts General Hospital, Boston, United States.
- ⁸⁵Medical University of South Carolina, Charleston, United States.
- ⁸⁶Department of Pediatrics, Division of Human Genetics, The Children's Hospital of Philadelphia and the Perelman School of Medicine at the University of Pennsylvania, Philadelphia, United States.
- ⁸⁷Center for Medical Information and Communication Technology, University Hospital Erlangen, Germany.
- ⁸⁸Renaissance Computing Institute/Department of Computer Science, University of North Carolina, Chapel Hill, Chapel Hill, United States.
- ⁸⁹Clinical Research Unit, Saint Antoine Hospital, APHP Greater Paris University Hospital, Paris, France.
- ⁹⁰Clevy.io, Paris, France.
- ⁹¹Department of Biostatistics, Harvard T. H. Chan School of Public Health, Boston, United States.
- ⁹²Department of Anaesthesia, National University Health Systems, Singapore, Singapore, Singapore.
- ⁹³Department of Biostatistics, Epidemiology and Informatics, University of Pennsylvania Perelman School of Medicine, Philadelphia, United States.
- ⁹⁴Harvard Catalyst, Harvard Medical School, Boston, United States.
- ⁹⁵Clinical Research Unit, Paris Saclay, APHP Greater Paris University Hospital, Boulogne-Billancourt, France.
- ⁹⁶Medical Informatics Center, Hospital das Clínicas, Faculty of Medicine of Botucatu, Clinical Research Unit of Botucatu Medical School, São Paulo State University, Botucatu, Brazil.
- ⁹⁷Department of Surgery, Beth Israel Deaconess Medical Center, Boston, United States.
- ⁹⁸Department of Anesthesiology and Critical Care, Children's Hospital of Philadelphia, Philadelphia, United States.
- ⁹⁹Department of Medical and Surgical Sciences, Infectious and Tropical Disease Unit, University Magna Graecia of Catanzaro, Catanzaro, Italy.
- ¹⁰⁰ENS, PSL University, Paris, France.
- ¹⁰¹Department of Medical and Surgical Sciences, University Magna Graecia of Catanzaro, Catanzaro, Italy.
- ¹⁰²Internal Medicine Department of Botucatu Medical School, São Paulo State University, Botucatu, Brazil.
- ¹⁰³Department of Biomedical Health Informatics, Children's Hospital of Philadelphia, Philadelphia, United States.
- ¹⁰⁴Department of Computational Biomedicine, Cedars-Sinai Medical Center, West Hollywood, United States.
- ¹⁰⁵Pediatric Infectious Disease Department, Hospital Universitario 12 de Octubre, Madrid, Spain.
- ¹⁰⁶Department of Biostatistics, Epidemiology, and Informatics, Institute for Biomedical Informatics, University of Pennsylvania Perelman School of Medicine, Berwyn, United States.
- ¹⁰⁷Department of Infectious Diseases, Great Ormond Street Hospital for Children, UK, London, United Kingdom.
- ¹⁰⁸Harvard Catalyst | The Harvard Clinical and Translational Science Center, Harvard Medical School, Boston, United States.
- ¹⁰⁹Department of Biomedical informatics, WiSDM, National University Health System Singapore, Singapore, Singapore.
- ¹¹⁰Department of Psychiatry, University of Pennsylvania Perelman School of Medicine, Philadelphia, United States.
- ¹¹¹Computational Health Informatics Program and Department of Pediatrics, Boston Children's Hospital, Harvard Medical School, Boston, United States.
- ¹¹²CTSI, WFBMI, Wake Forest School of Medicine, Winston Salem, United States.
- ¹¹³NC TraCS Institute, UNC Chapel Hill, Chapel Hill, United States.
- ¹¹⁴Department of Surgery, Beth Israel Deaconess Medical Center, Harvard Medical School, Boston, United States.
- ¹¹⁵Department of Medical Informatics, University of Erlangen-Nürnberg, Erlangen, Germany.
- ¹¹⁶Clinical Research Unit São Paulo State University, Brazil, Clinical Research Unit of Botucatu Medical School, São Paulo State University, Botucatu, Brazil.
- ¹¹⁷Office of Research and Development, Department of Veterans Affairs, Washington, DC, United States.
- ¹¹⁸Division of Infectious Diseases, Department of Medicine II, Medical Center – University of Freiburg, Faculty of Medicine, Freiburg, Germany.
- ¹¹⁹Pediatric Infectious Disease Department, Hospital Universitario 12 de Octubre, Madrid, Spain.
- ¹²⁰Biomedical Data Science Lab, ITACA Institute, Universitat Politècnica de València, Spain, Valencia, Spain.
- ¹²¹Department of Pediatrics (Critical Care), Northwestern University Feinberg School of Medicine, Chicago, United States.
- ¹²²Nurse department of FMB - medicine school of Botucatu, Clinical Research Unit of Botucatu Medical School, São Paulo State University, Botucatu, Brazil.
- ¹²³ASST Pavia, Lombardia Region Health System, Management Engineer, Direction, Pavia, Italy.
- ¹²⁴Data Analytics Center, University of Pennsylvania Health System, Philadelphia, United States.
- ¹²⁵Department of Anesthesiology, University Hospital Erlangen, FAU Erlangen-Nürnberg, Erlangen, Germany.
- ¹²⁶Hôpital Saint Louis, Department of Biostatistics and Bioinformatics, APHP Greater Paris University Hospital, Paris, France.
- ¹²⁷MICHR Informatics, University of Michigan, Ann Arbor, United States.

Harrison G Zhang^{1,2,5}, Jacqueline P Honerlaw², Monika Maripuri², Malarkodi Jebathilagam Samayamuthu³, Brendin R Beaulieu-Jones¹, Huma S Baig⁴, Sehi L'Yi¹, Yuk-Lam Ho², Michele Morris³, Vidul Ayakulangara Panickan¹, Xuan Wang¹, Griffin M Weber¹, Katherine P Liao^{2,5}, Shyam Visweswaran³, Bryce W.Q. Tan⁶, William Yuan¹, Nils Gehlenborg¹, Sumitra Muralidhar⁷, Rachel B Ramoni⁷, The Consortium for Clinical Characterization of COVID-19 by EHR (4CE)¹, Isaac S Kohane^{1,*}, Zongqi Xia^{8,*}, Kelly Cho^{2,*}, Tianxi Cai^{1,2,*}, Gabriel A Brat^{1,*†}

¹Department of Biomedical Informatics, Harvard Medical School, Boston, MA, United States.

²Massachusetts Veterans Epidemiology Research and Information Center (MAVERIC), VA Boston Healthcare System, Boston, MA, United States.

³Department of Biomedical Informatics, University of Pittsburgh, Pittsburgh, PA, United States.

⁴Department of Surgery, Beth Israel Deaconess Medical Center, Boston, MA, United States.

⁵Division of Rheumatology, Inflammation, and Immunity, Brigham and Women's Hospital, Boston, MA, United States.

⁶Department of Medicine, National University Hospital, Singapore, Singapore, Singapore.

⁷Office of Research and Development, U.S. Department of Veterans Affairs, Washington, D.C., United States.

⁸Department of Neurology, University of Pittsburgh, Pittsburgh, PA, United States

The value of large-scale real-world data such as that from electronic health records (EHRs) has been used to establish vaccine efficacy, elucidate the genetic etiologies of diseases, and advance epidemiological research.¹⁻³ Real-world data also has the potential to capture the wide spectrum of clinical features attributed to post-acute sequelae of SARS-CoV-2, also called long COVID, in diverse patient populations.⁴

We are an international consortium that has operationalized definitions of long COVID using health agency guidelines and established a chart review procedure based on these definitions.⁵ During this process, we identified 3 major challenges in using real world data to study long COVID: ambiguity and heterogeneity in clinical coding of long COVID; inadequacy of diagnostic codes in capturing the constellation of symptoms; and biases in

¹²⁸Department of Neurology, University of Pittsburgh Medical Center, Pittsburgh, United States.

¹²⁹Critical Care Medicine, Department of Medicine, St. Luke's University Health Network, Bethlehem, United States.

¹³⁰Department of Biomedical Health Informatics and the Department of Pediatrics, The Children's Hospital of Philadelphia and the University of Pennsylvania Perelman Medical School, Philadelphia, United States.

¹³¹Department of Information Management, National Central University, Taoyuan, Taiwan.

¹³²Clinical & Translational Science Institute, Medical College of Wisconsin, Milwaukee, United States.

¹³³Université Paris-Saclay, Inria, CEA, Montréal Neurological Institute, McGill University, Palaiseau, France.

¹³⁴SequeL, Inria Lille, Villeneuve-d'Ascq, France.

¹³⁵Respiratory Department, ICS S. Maugeri IRCCS Pavia Italy, Lumezzane (BS), ITALY.

¹³⁶Department of Health Management and Informatics, University of Missouri, Columbia, Columbia, United States.

¹³⁷Department of Anesthesiology and Critical Care Medicine, Children's Hospital of Philadelphia and University of Pennsylvania, Philadelphia, United States.

¹³⁸Department of Oncology, ASST Papa Giovanni XXIII, Bergamo, Bergamo, Italy

Competing interests

The authors declare no competing interests

EHR data arising from variability in the number and kind of contacts with the healthcare system. These challenges warrant special attention if the clinical community wishes to arrive at a robust understanding of long COVID using evidence derived from real world data.

We performed a manual medical record review of 300 randomly sampled patients infected with SARS-CoV-2 and assigned an International Classification of Diseases (ICD)-10 code (U09.9) for long COVID at the Beth Israel Deaconess Medical Center, University of Pittsburgh Medical Center, and national U.S. Veterans Health Administration.⁵ These three health systems collectively serve over 15 million patients each year.

We evaluated the extent to which patients with the ICD-10 code for this condition met our operationalized definitions from the World Health Organization's (WHO) and the US Centers for Disease Control.⁵⁻⁷ Our long COVID definition based on WHO guidelines required that a patient present with at least two new-onset persistent symptoms lasting for 60 days after infection, whereas our definition based on CDC guidelines required that a patient present with at least one new-onset persistent symptom lasting for 30 days.⁵

A comparison of real-world EHR and administrative data with manually extracted clinical information (obtained through chart review of patients with the U09.9 code) found that functional definitions of long COVID varied widely by provider, which led to inconsistencies in coding practice and adherence to clinical definitions. Among patients assigned the U09.9 code, an average of 40.2% met the more stringent WHO definition, 58.3% had a single symptom that met the WHO definition, and 65.4% met the least stringent CDC definition.⁵ This shows that the ICD-10 code is an unreliable surrogate of long COVID disease status in research. Research and policy efforts are needed to converge on a definition that will standardize coding practices and improve the ICD-10 code reliability.

Coding is further obfuscated by the potential for misclassification of long COVID with long-lasting complications from acute hospitalization, which are not specific to COVID-19.⁸ We found an average of 42.3% patients assigned the U09.9 code were hospitalized after infection and an average of 12.3% received intensive care, both of which can produce long lasting symptoms that overlap with long COVID.⁵ Physical and physiological effects of hospitalization or critical care are important patient-level factors that should not be misattributed to SARS-CoV-2 infection.

Capturing long COVID symptomology using diagnosis codes is difficult, as the syndrome encompasses a constellation of non-specific symptoms including pain, fatigue, and brain fog that are not well represented by coding schemes such as ICD-10.⁴⁻⁷ Leveraging textual data from EHRs may improve the ability to capture symptomatology. When we examined the data capture of symptoms by ICD-10 codes and natural language processing of clinical narratives, such as clinician notes and discharge summaries, we found that the incorporation of narrative data significantly improved identification of symptoms, compared to using diagnosis codes alone.⁵ This shows the potential use of natural language processing techniques to ascertain a more complete representation of a patient's health.

A further challenge is the definition of long COVID patient cohorts, as the syndrome is defined by a time to presentation and therefore requires an index date from which to observe

clinical outcomes. The index date is usually an initial infection date, which may become increasingly difficult to ascertain with the use of at-home testing, the results of which are inconsistently reported in EHRs. Researchers should therefore perform routine quality controls (such as checking the time period between initial infection and input of the ICD-10 code for long COVID), to better understand biases present in the data. Researchers should also allow for some flexibility in defining index dates, such as considering an infection time period rather than a single date, which helps to account for delays in billing or data processing.

Researchers must remain cognizant of potential patient selection bias in real-world data; using visits to a long COVID clinic as a proxy for true disease status is problematic. We found that, on average, only 24.0% of patients assigned the U09.9 code visited a long COVID clinic, suggesting that the majority of patients sampled were being coded by physicians who do not work at these clinics.⁵ Among patients who met the WHO definition of long COVID, only an average of 35.6% visited a long COVID clinic, suggesting that many patients are not being seen at these specialty care facilities.⁵

Studies should also account for differential data density and healthcare utilization. We found that patients who visited a long COVID clinic were on average annotated with more new-onset conditions when compared to patients who never visited a long COVID clinic.⁵ Physicians working at long COVID clinics could be more experienced with the syndrome and therefore document the disease more thoroughly. This contributes to a difference in data density and granularity, which can confound findings if not properly addressed.

Studies of long COVID using real-world data must be based on robust and comprehensive clinical datasets. The incorporation of narrative data obtained using natural language processing techniques should better capture symptoms, and researchers should take caution when using only the ICD-10 code or a visit to a long COVID clinic as surrogates for disease status.

Computational phenotypes (where data elements are combined using machine learning algorithms to describe a particular disorder) have the potential to account for the longitudinal persistence of symptoms while avoiding the misattribution of conditions that existed prior to initial infection.⁹ Semi-supervised machine learning algorithms are resistant to some of these challenges, and so may be powerful tools to capture complex underlying temporal patterns in the data using a small number of manually curated labels. Rule-based algorithms may be less suited for the inherent complexity of long COVID.¹⁰

Real-world data has an important role in supporting long COVID research, but these pitfalls should be considered so that the most equitable clinical and policy decisions can be informed by population-level studies.

REFERENCES

1. Haas EJ et al. *The Lancet* 397, 1819–1829 (2021).
2. McCarty CA et al. *BMC Med Genomics* 4, 13 (2011). [PubMed: 21269473]
3. Weber GM et al. *J Med Internet Res* 23, e31400 (2021). [PubMed: 34533459]

4. Zhang HG et al. NPJ Digit Med 5, 81 (2022). [PubMed: 35768548]
5. Zhang HG et al. 2023.02.12.23285701 Preprint at 10.1101/2023.02.12.23285701 (2023).
6. Soriano JB et al. Lancet Infect Dis 22, e102–e107 (2022). [PubMed: 34951953]
7. CDC. Post-COVID Conditions. Centers for Disease Control and Prevention <https://www.cdc.gov/coronavirus/2019-ncov/long-term-effects/index.html> (2022).
8. Clift AK et al. JAMA Psychiatry 79, 690–698 (2022). [PubMed: 35544272]
9. Liao KP et al. BMJ 350, h1885 (2015). [PubMed: 25911572]
10. Zhang Y et al. Nat Protoc 14, 3426–3444 (2019). [PubMed: 31748751]