

# Patient-Related Metadata Reported in Sequencing Studies of SARS-CoV-2: Protocol for a Scoping Review and Bibliometric Analysis

Karen O'Connor<sup>1</sup>, Davy Weissenbacher<sup>2</sup>, Amir Elyaderani<sup>3</sup>, Ebbing Lautenbach<sup>1,4,5</sup>, Matthew Scotch<sup>3,6</sup>, Graciela Gonzalez-Hernandez<sup>2</sup>

<sup>1</sup>Department of Biostatistics, Epidemiology and Informatics, Perelman School of Medicine, University of Pennsylvania, Philadelphia, PA, USA

<sup>2</sup>Department of Computational Biomedicine, Cedars-Sinai Medical Center, West Hollywood, CA, USA

<sup>3</sup>Biodesign Center for Environmental Health Engineering, Arizona State University, Tempe, AZ, USA

<sup>4</sup>Division of Infectious Diseases, Department of Medicine, Perelman School of Medicine, University of Pennsylvania, Philadelphia, Pennsylvania, USA.

<sup>5</sup>Center for Clinical Epidemiology and Biostatistics, Perelman School of Medicine, University of Pennsylvania, Philadelphia, Pennsylvania, USA.

<sup>6</sup>College of Health Solutions, Arizona State University, Tempe, AZ, USA

Corresponding Author:

Karen O'Connor

Senior Data Analyst

Department of Biostatistics, Epidemiology, and Informatics

Perelman School of Medicine

University of Pennsylvania

Email: [karoc@pennmedicine.upenn.edu](mailto:karoc@pennmedicine.upenn.edu)

## Abstract

**Background:** There has been an unprecedented effort to sequence the SARS-CoV-2 virus and examine its molecular evolution. This has been facilitated by the availability of publicly accessible databases, the Global Initiative on Sharing All Influenza Data (GISAID) and GenBank, which collectively hold millions of SARS-CoV-2 sequence records. Genomic epidemiology, however, seeks to go beyond phylogenetic analysis by linking genetic information to patient characteristics and disease outcomes, enabling a comprehensive understanding of transmission dynamics and disease impact.

While these repositories include fields reflecting patient-related metadata for a given sequence, inclusion of these demographic and clinical details is scarce. The extent to which patient-related metadata is reported in published sequencing studies and its quality remains largely unexplored.

**Methods:** The NIH's LitCovid collection will be used for automated classification of articles reporting having deposited SARS-CoV-2 sequences in public repositories, while an independent search will be conducted in PubMed for validation. Data extraction will be conducted using Covidence. The extracted data will be synthesized and summarized to quantify the availability of patient metadata in the published literature of SARS-CoV-2 sequencing studies. For the bibliometric analysis, relevant data points, such as author affiliations and citation metrics will be extracted.

**Discussion:** This scoping review will report on the extent and types of patient-related metadata reported in genomic viral sequencing studies of SARS-CoV-2, identify gaps in this reporting, and make recommendations for improving the quality and consistency of reporting in this area. The bibliometric analysis will uncover trends and patterns in the reporting of patient-related metadata, including differences in reporting based on study types or geographic regions. Co-occurrence networks of author keywords will also be presented. The insights gained from this study may help improve the quality and consistency of reporting patient metadata, enhancing the utility of sequence metadata and facilitating future research on infectious diseases.

**Keywords:** SARS-CoV-2; COVID-19; genomic epidemiology; scoping review; protocol

## Strengths and Limitations

- Classification of papers using machine learning will help pinpoint more relevant papers reducing time to sift through articles for inclusion.
- Study will provide an overview of the availability of information required to facilitate genomic epidemiology studies of SARS-CoV-2
- The bibliometric analysis will analyze publication trends, author affiliations, citation metrics, and other bibliographic information to provide several insights into the broader landscape of research that includes patient metadata reporting in genomic viral sequencing studies of SARS-CoV-2
- As with any search method for relevant studies, some may be missed by both the automated methods proposed and manually developed queries.

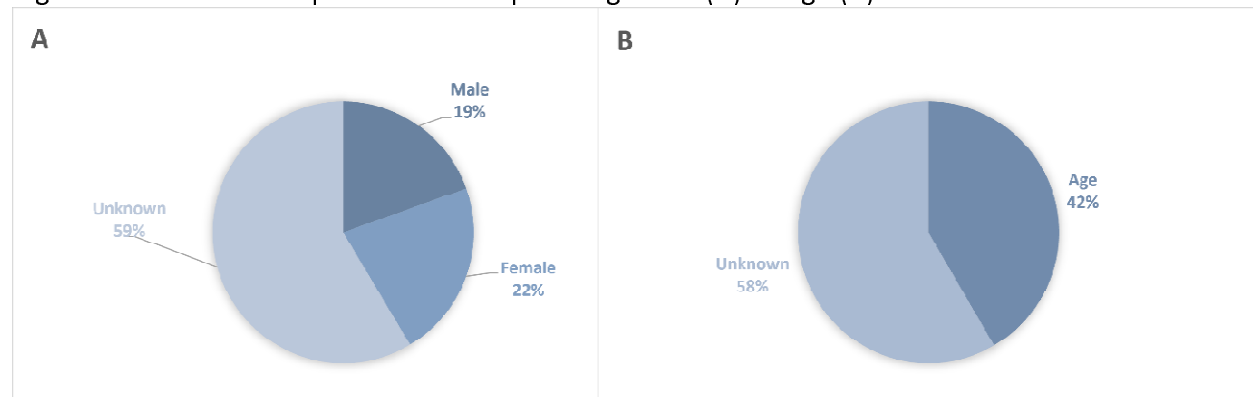
## Introduction

Since the onset of the COVID-19 pandemic, there has been an unprecedented effort in genomic epidemiology to sequence the virus, study its transmission, and examine molecular evolution. Public repositories, such as the Global Initiative on Sharing Avian Influenza Data (GISAID)<sup>1</sup> and NCBI's GenBank<sup>2</sup> host millions of SARS-CoV-2 sequence records. As of July 2023, GISAID contains 15.7 million sequences, while 7.7 million have been deposited in GenBank.

The availability of this vast amount of genomic data has facilitated significant discoveries, particularly in phylogenetic and phylodynamic studies<sup>3-5</sup>. Beyond phylogenetic studies, genomic epidemiology aims to understand the transmission dynamics, evolution, and impact of infectious diseases by analyzing the genetic information of pathogens and linking it to patient demographics and disease outcomes<sup>6,7</sup>. This work enables the tracking of the spread of pathogens, identifying high-risk populations, and discovering genetic factors that influence disease transmission, severity, and treatment response<sup>6,8</sup>. This knowledge can, in turn, inform public health strategies, guide the development of targeted interventions, and improve the overall understanding of infectious diseases<sup>9</sup>.

Ideally, patient geographic, demographic, and clinical information (such as disease severity and outcome) should be included in the sequence metadata upon its submission to the repository. Both GISAID and GenBank provide the location of the infected host (LOIH) information in their sequence metadata, however, the reported location granularity may vary and often lacks important details such as patient travel history. Similarly, patient demographic and clinical information is rarely complete: approximately 60% of sequences in GISAID have the age or gender of the infected host of the sequence entered as unknown (e.g., 'Not Available', 'Declined', 'Not Reported', etc.) (Figure 1), while GenBank lacks any standardized fields to include this information with sequence submissions.

Figure 1: Percent of sequences with reported gender (A) or age (B) in GISAID\*.



\*data downloaded April 3, 2023, representing 15.3 million sequences

This patient information, or at least a subset of it, may be reported in the published studies of those who obtained and performed the genomic sequencing. However, the extent to which location information as well as patient demographic or other clinical information is reported in

SARS-CoV-2 sequencing studies remains largely unexplored. Our review aims to bridge this gap in understanding by quantifying the extent and types of patient-related metadata reported in genomic viral sequencing studies of SARS-CoV-2.

Traditionally, identifying studies for a review requires the development of a detailed search strategy of databases using keywords and index terms, querying the title and abstracts of published articles. The selection of keywords greatly influences search results, leading to potential missed or irrelevant studies. Moreover, for the particular focus of our study, discussions of sequencing are often confined to the methods section of papers, rendering title and abstract screening less informative. While a substantial number of research articles related to SARS-CoV-2 and the pandemic have been published, there is sparse linkage between the sequence and publication databases making it difficult to identify publications relevant to the sequences. To overcome these limitations, we propose using an automated classifier to identify relevant studies for review.

A bibliometric analysis uses different methods and data points to quantify the trends and assess the impact of publications in a specific field<sup>10</sup>. While several bibliometric analyses have investigated COVID-19 related research trends, in general,<sup>11-13</sup> and in specific fields such as neurology<sup>14</sup>, long Covid<sup>15</sup> and medical imaging<sup>16</sup>, or specific geographic locations such as Africa<sup>17</sup>, no analysis has specifically focused on the publication trends related to the reporting of patient metadata related to SARS-CoV-2 genomic sequences. Our aims with this review and analysis are to identify reporting and publication trends as well as highlight the gaps in reporting that may hinder the advancement of genomic epidemiology studies of the COVID-19 pandemic.

### **Primary Research Objectives**

1. To quantitatively assess the extent and quality of patient-reported metadata, including demographic, clinical, and geographic information, in articles reporting original whole genome sequencing of the SARS-CoV-2 virus.
2. To perform a comprehensive bibliometric analysis to ascertain differences and discernible patterns between articles that include patient metadata and those that do not, thereby providing insights into the characteristics and factors associated with the reporting of patient data in the literature.
3. To evaluate the efficacy and reliability of a machine learning classifier in accurately identifying relevant articles for inclusion in the scoping review, enhancing the efficiency and effectiveness of the study selection process.

## **Methods**

Our review will follow the methodological framework identified by Arksey and O'Malley<sup>18</sup>. The scoping review will be reported in line with the PRISMA-SrC checklist<sup>19</sup>.

## Data Sources

We will utilize the NIH's LitCovid collection (1) for our machine learning classification. LitCovid is a curated collection of scholarly articles related to the Coronavirus Disease 2019 (COVID-19). The collection contains over 378,000 publications from 8,000 journals and is updated daily. LitCovid includes published articles as well as preprints. Additionally, we will independently search PubMed directly using a two-faceted search strategy and the NCBI E-utilities program to find publications linked to sequences. This combined approach will help ensure a comprehensive coverage of the literature for our study.

## Search Strategy

### Classification Model

Our classification model will be trained using manually annotated data. A full-text search strategy was developed to filter the LitCovid collection resulting in a corpus of targeted articles for annotation. The papers identified through the pipeline were annotated by two experienced annotators using the Inception annotation tool<sup>20</sup>. The annotators reviewed the full text of the articles and labeled sentences which confirmed the study's performance of SARS-CoV-2 sample sequencing from human specimens. The classifier will be instantiated as a pre-trained neural network, specifically a transformer model called BERT-base-uncased.

### Search Strategy

To evaluate our classifier and identify studies that may have been missed due to classification errors or lack of full text in the LitCovid collection, we will create a search strategy to independently search PubMed. We will develop a two-faceted search strategy to find "SARS-CoV-2" and "whole genome sequencing" related publications. We will utilize the search strategy developed for the LitCovid collection with additional keywords added to identify studies that report whole genome sequencing. Additionally, we will search for publications linked to SARS-CoV-2 sequences using the NCBI's E-utilities eLink programming API.

A publication date restriction of December 2019 onwards will be used in the searches as this review is focused on SARS-CoV-2 sequencing studies. No language restrictions will be placed on the searches, although financial and logistical restraints will not allow translation from all languages. All results will be uploaded to a Zotero<sup>21</sup> library where duplicate results will be removed.

## Inclusion/Exclusion Criteria

Papers positively identified by our classifier and our search results will be reviewed for inclusion in the review based on the criteria outlined in Table 1.

**Table 1. Inclusion and exclusion criteria for the scoping review.**

Facet	Inclusion Criteria	Exclusion Criteria
-------	--------------------	--------------------

<b>Sample Origin</b>	Individual human subject	<input type="checkbox"/> Non-human sources (e.g., mice, bats, ferrets) <input type="checkbox"/> Wastewater <input type="checkbox"/> Microbiome <input type="checkbox"/> Cloned/Cell Culture Virus
<b>Sequencing Type</b>	Whole genomic sequencing	Studies will be excluded if the following sequencing methods were exclusively performed: <ul style="list-style-type: none"> <li><input type="checkbox"/> PCR or LAMP for viral detection</li> <li><input type="checkbox"/> Single-cell sequencing</li> <li><input type="checkbox"/> Gene expression studies</li> <li><input type="checkbox"/> Protocol validation studies on cell culture virus</li> <li><input type="checkbox"/> Exome sequencing</li> </ul>
<b>Study Design</b>	Any type of peer-reviewed or preprint study reporting on the original sequencing of SARS-CoV-2 samples.	Any other study design, or a study not peer reviewed.
<b>Publication Dates</b>	December 2019, or later	Before December 2019
<b>Language</b>	All	None

Two reviewers will perform title and abstract screening using the Covidence<sup>22</sup> systematic review management tool with any disagreements resolved by discussion. Two independent reviewers will also conduct full-text screening in Covidence.

## Data Extraction

Data extraction will be conducted in Covidence. The reviewers will examine the full text of the articles, including any supplementary files, for data extraction. The customizable interface will be designed to prompt the reviewer to extract various details, such as general publication information, study characteristics, sequencing specifics, and the presence or absence of patient demographic, clinical, or location information. Furthermore, the location of this information within the articles will be noted. An example of the data extraction form can be found in Table 2.

**Table 2: Example of data that will be extracted from included studies.**

Prompt	Response
Publication Information	

Study Name	Free text
Article Title	Free text
Year of Publication	YYYY
Publication Type	Journal, Conference, Preprint
<b>Study and Sequence Information</b>	
Study Objective	Free text
Location of Study (country)	Free text
Number of patients	Free text
Number of samples sequenced	Free text
Repository sequences deposited to	GISAID, GenBank, Other, NR
For studies with >1 patient, are sequences linked to a patient?	Yes   No
<b>Patient Demographic Information Reported</b>	
Age	Yes   No
Gender	Yes   No
Race/Ethnicity	Yes   No
<b>Patient Clinical Information Reported</b>	
Symptoms	Yes   No
Severity	Yes   No
Inpatient or Outpatient	Yes   No
Treatments	Yes   No
Outcomes	Yes   No
<b>Patient Geographic Information Reported</b>	
Location of Residence	Yes   No
Travel Information	Yes   No

We will test the initial extraction form on a subset of articles and revise it as needed.

For bibliometric analysis, all pertinent data points will be extracted for studies included in our review including, author location and institution information, journal, study type, citation metrics, and author keywords.

## Data Analysis

The extracted data will be synthesized and summarized to quantify the availability of patient metadata in the published literature of SARS-CoV-2 sequencing studies using an exported spreadsheet from Covidence. For the bibliometric analysis, data will be analyzed and visualized using the VOSviewer<sup>23</sup> software or the bibliometrix<sup>24</sup> package for R. Other software tools may be used as needed for analysis.

## Results

We will summarize and narratively describe our findings, using tables, graphs, and charts when applicable regarding the number of sequences covered in our included studies, the distribution of the sequences in the respective repositories, and the quantity and type of reported patient metadata in the studies. We will also present the geographical location of the study's authors using maps and report our findings, including the most frequent journals and article types, as well as analyze differences between studies that reported patient data from those that did not. Co-occurrence networks of author keywords will be presented to highlight the frequency and differences in themes and study focus between the reporting groups.

## Discussion

There has been an unprecedented effort in the sequencing and sharing of the viral genomes of SARS-CoV-2 through publicly available databases, such as GISAID and GenBank, during the COVID-19 pandemic. However, the utility of these sequences for genomic epidemiology may not be fully realized due to the unavailability of relevant metadata about the patient from whom the specimen was obtained<sup>25</sup>. Our study aims to conduct a scoping review and bibliometric analysis focusing on patient metadata reporting in genomic viral sequencing studies of SARS-CoV-2.

This scoping review will provide valuable insights into the current state of reporting of patient-related metadata in SARS-CoV-2 sequencing studies. The review findings will be used to identify gaps in the reporting of patient metadata and make recommendations for improving the quality and consistency of reporting of patient-related metadata in SARS-CoV-2 sequencing studies.

In addition to the findings of our scoping review, the bibliometric analysis will likely identify several other important trends and patterns in the reporting of patient-related metadata. For example, the analysis may find that the reporting of patient-related metadata is more common in certain types of studies, or that it is more likely to be reported in studies from certain geographic regions. The findings of the scoping review and bibliometric analysis will provide valuable insights into the factors that influence the reporting of patient-related metadata and will help to inform future research on this topic. Furthermore, the identification and quantification of the metadata in literature may aid in advancing other research, such as the development of machine learning methods to extract this information and enhance sequence data through automatic methods.

## Strengths/Limitations

This study will conduct a systematic and comprehensive scoping review, encompassing a large number of articles from various databases ensuring a thorough examination of the current state of reporting patient-related metadata in SARS-CoV-2 sequencing studies, and will provide a comprehensive overview of the available literature. The inclusion of bibliometric analysis will



go beyond the scoping review to analyze publication trends, author affiliations, citation metrics, and other bibliographic information to provide several insights into the broader landscape of research that includes patient metadata reporting in genomic viral sequencing studies of SARS-CoV-2.

While some relevant studies may be missed due to search limitations and the classification model, our study's strength lies in providing valuable insights into the current state of reporting patient-related metadata. Although certain limitations exist, such as potential limitations in reported patient metadata<sup>26,27</sup> and the focus on SARS-CoV-2 sequencing studies, our findings will contribute to improving the quality and consistency of reporting in genomic epidemiology. Future research can build upon our study to address these gaps and enhance reporting practices in this field.

## Conclusion

This protocol outlines the steps that we will take in our scoping review which will be supported by an automated classifier and bibliometric analysis. We will fill the knowledge gap regarding the extent and types of patient-related metadata reported in genomic viral sequencing studies of SARS-CoV-2 and will provide valuable insights by identifying themes and trends in the published literature. The results of this study may encourage improved and standardized reporting practices which will significantly enhance the utility of sequence metadata and aid in advancing our understanding of the SARS-CoV-2 or any future pandemic.

## Declarations

### Ethics and Dissemination

This study will analyze and synthesize previously published information. Data sharing is not applicable to this article as no datasets were generated or analysed during the current study. We will submit for publication the completed scoping review and bibliometric analysis. At that time, any extracted data and data used in our analysis will be made available with the publication.

### Competing interests

The authors declare that they have no competing interests.

### Funding

Research reported in this publication was supported by the National Institute of Allergy And Infectious Diseases of the National Institutes of Health under Award Number R01AI164481 to GGH and MS. The NIH National Institute of Allergy And Infectious Diseases funded this research but were not involved in the conceptualization, design, data collection, analysis, decision to publish, or preparation of the manuscript. The views expressed in this manuscript are those of the authors and not those of the NIH.

### Authors' contributions

KO, EL, MS and GGH designed the study. KO was a major contributor in the writing of the manuscript. DW designed the classification methods. KO and AE designed the annotation methods. All authors read, edited and approved the final manuscript.

### Acknowledgements

Not applicable.

## References

1. Shu Y, McCauley J. GISAID: Global initiative on sharing all influenza data – from vision to reality. *Eurosurveillance*. 2017;22(13):1-1. doi:10.2807/1560-7917.ES.2017.22.13.30494
2. Sayers EW, Cavanaugh M, Clark K, Ostell J, Pruitt KD, Karsch-Mizrachi I. GenBank. *Nucleic Acids Res*. 2019;47(D1):D94-D99. doi:10.1093/nar/gky989
3. Forster P, Forster L, Renfrew C, Forster M. Phylogenetic network analysis of SARS-CoV-2 genomes. *Proc Natl Acad Sci*. 2020;117(17):9241-9243. doi:10.1073/pnas.2004999117
4. van Dorp L, Acman M, Richard D, et al. Emergence of genomic diversity and recurrent mutations in SARS-CoV-2. *Infect Genet Evol*. 2020;83:104351. doi:10.1016/j.meegid.2020.104351
5. Tang X, Wu C, Li X, et al. On the origin and continuing evolution of SARS-CoV-2. *Natl Sci Rev*. 2020;7(6):1012-1023. doi:10.1093/nsr/nwaa036
6. Hill V, Ruis C, Bajaj S, Pybus OG, Kraemer MUG. Progress and challenges in virus genomic epidemiology. *Trends Parasitol*. 2021;37(12):1038-1049. doi:10.1016/j.pt.2021.08.007
7. Tang P, Croxen MA, Hasan MR, Hsiao WWL, Hoang LM. Infection control in the new age of genomic epidemiology. *Am J Infect Control*. 2017;45(2):170-179. doi:10.1016/j.ajic.2016.05.015
8. National Academies of Sciences, Engineering, and Medicine. *Genomic Epidemiology Data Infrastructure Needs for SARS-CoV-2: Modernizing Pandemic Response Strategies*. The National Academies Press; 2020.
9. World Health Organization. *Genomic Sequencing of SARS-CoV-2 A Guide to Implementation for Maximum Impact on Public Health.*; 2021.
10. Gutiérrez-Salcedo M, Martínez MÁ, Moral-Munoz JA, Herrera-Viedma E, Cobo MJ. Some bibliometric procedures for analyzing and evaluating research fields. *Appl Intell*. 2018;48(5):1275-1287. doi:10.1007/s10489-017-1105-y
11. Hossain MM. Current Status of Global Research on Novel Coronavirus Disease (COVID-19): A Bibliometric Analysis and Knowledge Mapping. Published online May 18, 2020. doi:10.2139/ssrn.3547824
12. Nasab FR, Rahim F. *Bibliometric Analysis of Global Scientific Research on SARS-CoV-2 (COVID-19)*. Health Informatics; 2020. doi:10.1101/2020.03.19.20038752
13. Yu Y, Li Y, Zhang Z, et al. A bibliometric analysis using VOSviewer of publications on COVID-19. *Ann Transl Med*. 2020;8(13):816-816. doi:10.21037/atm-20-4235

14. Zhang Q, Li J, Weng L. A bibliometric analysis of COVID-19 publications in neurology by using the visual mapping method. *Front Public Health*. 2022;10. Accessed July 11, 2023. <https://www.frontiersin.org/articles/10.3389/fpubh.2022.937008>
15. Kim TH, Jeon SR, Kang JW, Kwon S. Complementary and Alternative Medicine for Long COVID: Scoping Review and Bibliometric Analysis. *Evid-Based Complement Altern Med ECAM*. 2022;2022:7303393. doi:10.1155/2022/7303393
16. Wen R, Zhang M, Xu R, et al. COVID-19 imaging, where do we go from here? Bibliometric analysis of medical imaging in COVID-19. *Eur Radiol*. 2023;33(5):3133-3143. doi:10.1007/s00330-023-09498-z
17. Guleid FH, Oyando R, Kabia E, Mumbi A, Akech S, Barasa E. A bibliometric analysis of COVID-19 research in Africa. *BMJ Glob Health*. 2021;6(5):e005690. doi:10.1136/bmjgh-2021-005690
18. Arksey H, O'Malley L. Scoping studies: towards a methodological framework. *Int J Soc Res Methodol*. 2005;8(1):19-32. doi:10.1080/1364557032000119616
19. Tricco AC, Lillie E, Zarin W, et al. PRISMA Extension for Scoping Reviews (PRISMA-ScR): Checklist and Explanation. *Ann Intern Med*. 2018;169(7):467-473. doi:10.7326/M18-0850
20. Klie JC, Bugert M, Boullosa B, de Castilho RE, Gurevych I. The INCEption Platform: Machine-Assisted and Knowledge-Oriented Interactive Annotation.
21. Zotero | Your personal research assistant. Accessed July 14, 2023. <https://www.zotero.org/>
22. Covidence - Better systematic review management. Covidence. Accessed July 14, 2023. <https://www.covidence.org/>
23. VOSviewer - Visualizing scientific landscapes. VOSviewer. Accessed July 11, 2023. <https://www.vosviewer.com//>
24. Aria M, Cuccurullo C. bibliometrix: An R-tool for comprehensive science mapping analysis. *J Informetr*. 2017;11(4):959-975. doi:10.1016/j.joi.2017.08.007
25. Grad YH, Lipsitch M. Epidemiologic data and pathogen genome sequences: a powerful synergy for public health. *Genome Biol*. 2014;15(11):538. doi:10.1186/s13059-014-0538-4
26. Hernandez MM, Gonzalez-Reiche AS, Alshammary H, et al. Molecular evidence of SARS-CoV-2 in New York before the first pandemic wave. *Nat Commun*. 2021;12(1):3463. doi:10.1038/s41467-021-23688-7
27. Page AJ, Mather AE, Le-Viet T, et al. Large-scale sequencing of SARS-CoV-2 genomes from one region allows detailed epidemiology and enables local outbreak management. *Microb Genomics*. 2021;7(6):000589. doi:10.1099/mgen.0.000589

