


# Tracking the international spread of SARS-CoV-2 lineages B.1.1.7 and B.1.351/501Y-V2 with grinch

## Other Journal Item

### Author(s):

O'Toole, Áine; Hill, Verity; Pybus, Oliver G.; Watts, Alexander; Bogoch, Issac I.; Khan, Kamran; Messina, Jane P.; Tegally, Houriiyah; Lessells, Richard R.; Giandhari, Jennifer; Pillay, Sureshnee; Tumed, Kefentse Arnold; Nyepetsi, Gape; Kebabonye, Malebogo; Matsheka, Maitshwarelo; Mine, Madisa; Tokajian, Sima; Hassan, Hamad; Salloum, Tamara; Merhi, Georgi; [Stadler, Tanja](#)  et al.

### Publication date:

2021-05-19

### Permanent link:

<https://doi.org/10.3929/ethz-b-000596303>

### Rights / license:

[Creative Commons Attribution 4.0 International](#)

### Originally published in:

Wellcome Open Research 6, <https://doi.org/10.12688/wellcomeopenres.16661.2>




SOFTWARE TOOL ARTICLE

**REVISED** Tracking the international spread of SARS-CoV-2**lineages B.1.1.7 and B.1.351/501Y-V2 with grinch [version 2; peer review: 3 approved]**

Previously titled: Tracking the international spread of SARS-CoV-2 lineages B.1.1.7 and B.1.351/501Y-V2

Áine O'Toole <sup>1\*</sup>, Verity Hill<sup>1\*</sup>, Oliver G. Pybus<sup>2</sup>, Alexander Watts<sup>3,4</sup>,  
Issac I. Bogoch <sup>5,6</sup>, Kamran Khan<sup>3-5</sup>, Jane P. Messina<sup>7</sup>,  
The COVID-19 Genomics UK (COG-UK) consortium,  
Network for Genomic Surveillance in South Africa (NGS-SA),  
Brazil-UK CADDE Genomic Network, Houriiyah Tegally<sup>8</sup>, Richard R. Lessells <sup>8</sup>,  
Jennifer Giandhari<sup>8</sup>, Sureshnee Pillay <sup>8</sup>, Kefentse Arnold Tumedi <sup>9</sup>,  
Gape Nyepetsi<sup>10</sup>, Malebogo Keabonye<sup>11</sup>, Maitshwarelo Matsheka <sup>9</sup>,  
Madisa Mine<sup>10</sup>, Sima Tokajian <sup>12</sup>, Hamad Hassan<sup>13</sup>, Tamara Salloum <sup>12</sup>,  
Georgi Merhi <sup>12</sup>, Jad Koweyes <sup>12</sup>, Jemma L. Geoghegan<sup>14,15</sup>, Joep de Ligt <sup>15</sup>,  
Xiaoyun Ren<sup>15</sup>, Matthew Storey<sup>15</sup>, Nikki E. Freed<sup>16</sup>, Chitra Pattabiraman <sup>17</sup>,  
Pramada Prasad<sup>17</sup>, Anita S. Desai<sup>17</sup>, Ravi Vasanthapuram<sup>17</sup>, Thomas F. Schulz <sup>18</sup>,  
Lars Steinbrück <sup>18</sup>, Tanja Stadler<sup>19</sup>, Swiss Viollier Sequencing Consortium,  
Antonio Parisi<sup>20</sup>, Angelica Bianco<sup>20</sup>, Darío García de Viedma<sup>21,22</sup>,  
Sergio Buenestado-Serrano <sup>21</sup>, Vítor Borges <sup>23</sup>, Joana Isidro<sup>23</sup>, Sílvia Duarte<sup>24</sup>,  
João Paulo Gomes<sup>23</sup>, Neta S. Zuckerman<sup>25</sup>, Michal Mandelboim<sup>25</sup>, Orna Mor<sup>25</sup>,  
Torsten Seemann<sup>26</sup>, Alicia Arnott<sup>27</sup>, Jenny Draper<sup>27</sup>, Mailie Gall <sup>27</sup>,  
William Rawlinson<sup>28</sup>, Ira Deveson <sup>29</sup>, Sanmarié Schlebusch<sup>30</sup>, Jamie McMahon<sup>30</sup>,  
Lex Leong<sup>31</sup>, Chuan Kok Lim<sup>31</sup>, Maria Chironna<sup>32</sup>, Daniela Loconsole <sup>32</sup>,  
Antonin Bal<sup>33</sup>, Laurence Josset<sup>33</sup>, Edward Holmes<sup>34</sup>, Kirsten St. George<sup>35</sup>,  
Erica Lasek-Nesselquist<sup>35</sup>, Reina S. Sikkema<sup>36</sup>, Bas Oude Munnink<sup>36</sup>,  
Marion Koopmans<sup>36</sup>, Mia Brytting<sup>37</sup>, V. Sudha rani <sup>38</sup>, S. Pavani<sup>38</sup>,  
Teemu Smura<sup>39</sup>, Albert Heim<sup>18</sup>, Satu Kurkela<sup>40</sup>, Massab Umair<sup>41</sup>,  
Muhammad Salman<sup>41</sup>, Barbara Bartolini<sup>42</sup>, Martina Rueca<sup>42</sup>, Christian Drosten<sup>43</sup>,  
Thorsten Wolff <sup>44</sup>, Olin Silander<sup>16</sup>, Dirk Eggink<sup>45</sup>, Chantal Reusken<sup>45</sup>,  
Harry Vennema<sup>45</sup>, Aekyung Park<sup>46</sup>, Christine Carrington <sup>47</sup>, Nikita Sahadeo<sup>47</sup>,  
Michael Carr<sup>48</sup>, Gabo Gonzalez<sup>48</sup>, SEARCH Alliance San Diego,  
National Virus Reference Laboratory, SeqCOVID-Spain,

Danish Covid-19 Genome Consortium (DCGC),  
Communicable Diseases Genomic Network (CDGN),  
Dutch National SARS-CoV-2 surveillance program,  
Division of Emerging Infectious Diseases (KDCA), Tulio de Oliveira<sup>8</sup>,  
Nuno Faria <sup>2,49</sup>, Andrew Rambaut<sup>1</sup>, Moritz U. G. Kraemer<sup>2</sup>

<sup>1</sup>Institute of Evolutionary Biology, University of Edinburgh, Edinburgh, UK

<sup>2</sup>Department of Zoology, University of Oxford, Oxford, UK

<sup>3</sup>Li Ka Shing Knowledge Institute, St. Michael's Hospital, Toronto, Canada

<sup>4</sup>BlueDot, Toronto, Canada

<sup>5</sup>Department of Medicine, University of Toronto, Toronto, Canada

<sup>6</sup>Divisions of General Internal Medicine and Infectious Diseases, University Health Network, Toronto, Canada

<sup>7</sup>Department of Geography, University of Oxford, Oxford, UK

<sup>8</sup>KwaZulu-Natal Research Innovation and Sequencing Platform (KRISP), Nelson R Mandela School of Medicine, University of KwaZulu-Natal, Durban, South Africa

<sup>9</sup>Botswana Institute for Technology Research and Innovation, Gaborone, Botswana

<sup>10</sup>National Health Laboratory, Gaborone, Botswana

<sup>11</sup>Ministry of Health and Wellness, Gaborone, Botswana

<sup>12</sup>Department of Natural Sciences, Lebanese American University, Beirut, Lebanon

<sup>13</sup>Faculty of Public Health, Lebanese University, Beirut, Lebanon

<sup>14</sup>Department of Microbiology and Immunology, University of Otago, Dunedin, New Zealand

<sup>15</sup>Institute of Environmental Science and Research, Wellington, New Zealand

<sup>16</sup>School of Natural and Computational Sciences, Massey University, Auckland, New Zealand

<sup>17</sup>Department of Neurovirology, National Institute of Mental Health and Neurosciences, Bengaluru, India

<sup>18</sup>Institute of Virology, Hannover Medical School, Hannover, Germany

<sup>19</sup>Department of Biosystems Science and Engineering, ETH Zürich, Zurich, Switzerland

<sup>20</sup>Istituto Zooprofilattico sperimentale della Puglia e della Basilicata, Puglia, Italy

<sup>21</sup>Hospital General Universitario Gregorio Marañón; Instituto de Investigación Sanitaria Gregorio Marañón, Madrid, Spain

<sup>22</sup>CIBER Enfermedades Respiratorias CIBERES, Madrid, Spain

<sup>23</sup>Bioinformatics Unit, Department of Infectious Diseases, National Institute of Health Doutor Ricardo Jorge (INSA), Lisbon, Portugal

<sup>24</sup>Innovation and Technology Unit, Department of Human Genetics, National Institute of Health Doutor Ricardo Jorge (INSA), Lisbon, Portugal

<sup>25</sup>Central Virology Laboratory, Israel Ministry of Health, Sheba Medical Center, Ramat Gan, Israel

<sup>26</sup>Microbiological Diagnostic Unit Public Health Laboratory, Department of Microbiology & Immunology, University of Melbourne at the Peter Doherty Institute for Infection & Immunity, Melbourne, Australia

<sup>27</sup>New South Wales Health Pathology - Institute of Clinical Pathology and Medical Research, Sydney, Australia

<sup>28</sup>New South Wales Health Pathology Randwick, Prince of Wales Hospital, Sydney, Australia

<sup>29</sup>Kinghorn Centre for Clinical Genomics, Sydney, Australia

<sup>30</sup>Queensland Reference Centre for Microbial and Public Health Genomics, Forensic and Scientific Services, Health Support Queensland, Queensland Health South Australia Pathology, Adelaide, Australia

<sup>31</sup>South Australia Pathology, Adelaide, Australia

<sup>32</sup>Department of Biomedical Sciences and Human Oncology, University of Bari, Bari, Italy

<sup>33</sup>Centre National de Référence des virus des infections respiratoires, Hospices Civils de Lyon, Lyon, France

<sup>34</sup>University of Sydney, Sydney, Australia

<sup>35</sup>Wadsworth Center, New York State Department of Health, Albany, New York, USA

<sup>36</sup>ErasmusMC, Department of Viroscience, WHO collaborating centre for arbovirus and viral hemorrhagic fever Reference and Research, Rotterdam, The Netherlands

<sup>37</sup>The Public Health Agency of Sweden, Department of Microbiology, Solna, Sweden

<sup>38</sup>Upgraded Department of Microbiology, Osmania Medical College, Hyderabad, Telangana, India

<sup>39</sup>Department of Virology, University of Helsinki, Helsinki, Finland

<sup>40</sup>HUS Diagnostic Center, HUSLAB, Clinical Microbiology, University of Helsinki and Helsinki University Hospital, Helsinki, Finland

<sup>41</sup>Department of Virology, National Institute of Health, Islamabad, Pakistan

<sup>42</sup>National Institute for Infectious Diseases "L. Spallanzani", Rome, Italy

<sup>43</sup>Institute for Virology, Charité Universitätsmedizin, Berlin, Germany

<sup>44</sup>Robert Koch-Institut, , Head, Unit 17, Influenza and other Respiratory Viruses, Seestr. 10, Berlin, Germany

<sup>45</sup>WHO COVID-19 reference laboratory, Centre for Infectious Disease Control-National Institute for Public Health and the Environment, Bilthoven, The Netherlands

<sup>46</sup>Division of Emerging Infectious Diseases, Bureau of Infectious Disease Diagnosis Control, Korea Disease Control and Prevention Agency, Cheongju-si, Chungcheongbuk-do, South Korea

<sup>47</sup>University of the West Indies, St. Augustine, Trinidad and Tobago

<sup>48</sup>National Virus Reference Laboratory, University College Dublin, Dublin, Ireland

<sup>49</sup>Imperial College London, London, UK

\* Equal contributors

**v2** First published: 19 May 2021, 6:121  
<https://doi.org/10.12688/wellcomeopenres.16661.1>

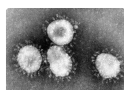
Latest published: 17 Sep 2021, 6:121  
<https://doi.org/10.12688/wellcomeopenres.16661.2>

## Abstract

Late in 2020, two genetically-distinct clusters of severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2) with mutations of biological concern were reported, one in the United Kingdom and one in South Africa. Using a combination of data from routine surveillance, genomic sequencing and international travel we track the international dispersal of lineages B.1.1.7 and B.1.351 (variant 501Y-V2). We account for potential biases in genomic surveillance efforts by including passenger volumes from location of where the lineage was first reported, London and South Africa respectively. Using the software tool grinch (global report investigating novel coronavirus haplotypes), we track the international spread of lineages of concern with automated daily reports. Further, we have built a custom tracking website ([cov-lineages.org/global\\_report.html](http://cov-lineages.org/global_report.html)) which hosts this daily report and will continue to include novel SARS-CoV-2 lineages of concern as they are detected.

## Keywords

genomic surveillance, air travel, SARS-CoV-2, genomics, genome sequencing, virus, surveillance, pandemic, B.1.1.7, B.1.351, N501Y, coronavirus, sequencing, genomic epidemiology



This article is included in the [Coronavirus \(COVID-19\)](#) collection.

## Open Peer Review

Approval Status ✓✓✓

	1	2	3
<b>version 2</b> (revision) 17 Sep 2021			
<b>version 1</b> 19 May 2021	✓ <a href="#">view</a>	✓ <a href="#">view</a>	✓ <a href="#">view</a>

- Rob Lanfear** , Australian National University, Canberra, Australia
- Anderson F. Brito** , Yale School of Public Health, New Haven, USA
- George Githinji** , Kenya Medical Research Institute (KEMRI) - Wellcome Trust Research Programme, Kilifi, Kenya

Any reports and responses or comments on the article can be found at the end of the article.

**Corresponding authors:** Áine O'Toole ([aine.otoole@ed.ac.uk](mailto:aine.otoole@ed.ac.uk)), Moritz U. G. Kraemer ([moritz.kraemer@zoo.ox.ac.uk](mailto:moritz.kraemer@zoo.ox.ac.uk))

**Author roles:** **O'Toole Á:** Conceptualization, Data Curation, Formal Analysis, Funding Acquisition, Investigation, Methodology, Project Administration, Software, Validation, Visualization, Writing – Original Draft Preparation, Writing – Review & Editing; **Hill V:** Conceptualization, Data Curation, Formal Analysis, Methodology, Software, Visualization, Writing – Original Draft Preparation, Writing – Review & Editing; **Pybus OG:** Conceptualization, Supervision, Writing – Original Draft Preparation, Writing – Review & Editing; **Watts A:** Data Curation, Resources, Visualization; **Bogoch II:** Data Curation, Resources, Visualization; **Khan K:** Data Curation, Resources, Visualization; **Messina JP:** Data Curation, Resources; **Tegally H:** Data Curation, Formal Analysis, Resources, Writing – Review & Editing; **Lessells RR:** Data Curation, Resources; **Giandhari J:** Data Curation, Formal Analysis, Resources; **Pillay S:** Data Curation, Formal Analysis,

Resources; **Tumedi KA**: Data Curation, Formal Analysis, Resources; **Nyepetsi G**: Data Curation, Resources; **Kebabonye M**: Resources; **Matsheka M**: Resources, Supervision; **Mine M**: Resources, Supervision; **Tokajian S**: Data Curation, Formal Analysis, Resources; **Hassan H**: Data Curation, Formal Analysis, Resources; **Salloum T**: Data Curation, Formal Analysis, Resources; **Merhi G**: Data Curation, Formal Analysis, Resources; **Koweyes J**: Data Curation, Formal Analysis, Resources; **Geoghegan JL**: Data Curation, Formal Analysis, Resources; **de Ligt J**: Data Curation, Formal Analysis, Resources; **Ren X**: Data Curation, Formal Analysis, Resources; **Storey M**: Data Curation, Formal Analysis, Resources; **Freed NE**: Data Curation, Formal Analysis, Resources; **Pattabiraman C**: Data Curation, Formal Analysis, Resources; **Prasad P**: Data Curation, Formal Analysis, Resources; **Desai AS**: Data Curation, Formal Analysis, Resources; **Vasanthapuram R**: Data Curation, Formal Analysis, Resources; **Schulz TF**: Data Curation, Formal Analysis, Resources; **Steinbrück L**: Data Curation, Formal Analysis, Resources; **Stadler T**: Data Curation, Formal Analysis, Resources, Supervision; **Parisi A**: Data Curation, Formal Analysis, Resources; **Bianco A**: Data Curation, Formal Analysis, Resources; **García de Viedma D**: Data Curation, Formal Analysis, Resources; **Buenestado-Serrano S**: Data Curation, Formal Analysis, Resources; **Borges V**: Data Curation, Formal Analysis, Resources; **Isidro J**: Data Curation, Formal Analysis, Resources; **Duarte S**: Data Curation, Formal Analysis, Resources; **Gomes JP**: Data Curation, Formal Analysis, Resources; **Zuckerman NS**: Data Curation, Formal Analysis, Resources; **Mandelboim M**: Data Curation, Formal Analysis, Resources; **Mor O**: Data Curation, Resources; **Seemann T**: Data Curation, Formal Analysis, Resources; **Arnott A**: Data Curation, Formal Analysis, Resources; **Draper J**: Data Curation, Formal Analysis, Resources; **Gall M**: Data Curation, Formal Analysis, Resources; **Rawlinson W**: Data Curation, Formal Analysis, Resources; **Deveson I**: Data Curation, Formal Analysis, Resources; **Schlebusch S**: Data Curation, Formal Analysis, Resources; **McMahon J**: Data Curation, Resources; **Leong L**: Data Curation, Formal Analysis, Resources; **Lim CK**: Data Curation, Formal Analysis, Resources; **Chironna M**: Data Curation, Formal Analysis, Resources; **Loconsole D**: Data Curation, Formal Analysis, Resources; **Bal A**: Data Curation, Formal Analysis, Resources; **Josset L**: Data Curation, Formal Analysis, Resources; **Holmes E**: Investigation, Writing – Original Draft Preparation, Writing – Review & Editing; **St. George K**: Data Curation, Formal Analysis, Resources; **Lasek-Nesselquist E**: Data Curation, Formal Analysis, Resources; **Sikkema RS**: Data Curation, Formal Analysis, Resources; **Oude Munnink B**: Data Curation, Formal Analysis, Resources; **Koopmans M**: Data Curation, Formal Analysis, Resources; **Brytting M**: Data Curation, Formal Analysis, Resources; **Sudha rani V**: Data Curation, Resources; **Pavani S**: Data Curation, Resources; **Smura T**: Data Curation, Formal Analysis, Resources; **Heim A**: Data Curation, Formal Analysis, Resources; **Kurkela S**: Data Curation, Formal Analysis, Resources; **Umair M**: Data Curation, Formal Analysis, Resources; **Salman M**: Data Curation, Formal Analysis, Resources; **Bartolini B**: Data Curation, Formal Analysis, Resources; **Rueca M**: Data Curation, Formal Analysis, Resources; **Drosten C**: Data Curation, Formal Analysis, Resources; **Wolff T**: Data Curation, Formal Analysis, Resources; **Silander O**: Data Curation, Formal Analysis, Resources; **Eggink D**: Data Curation, Resources; **Reusken C**: Data Curation, Resources; **Vennema H**: Data Curation, Formal Analysis, Resources; **Park A**: Data Curation, Formal Analysis, Resources; **Carrington C**: Data Curation, Formal Analysis, Resources; **Sahadeo N**: Data Curation, Formal Analysis, Resources; **Carr M**: Data Curation, Formal Analysis, Resources; **Gonzalez G**: Data Curation, Formal Analysis, Resources; **de Oliveira T**: Data Curation, Formal Analysis, Resources; **Faria N**: Data Curation, Formal Analysis, Investigation, Resources; **Rambaut A**: Conceptualization, Formal Analysis, Funding Acquisition, Supervision, Validation, Writing – Original Draft Preparation, Writing – Review & Editing; **Kraemer MUG**: Conceptualization, Data Curation, Formal Analysis, Investigation, Project Administration, Resources, Supervision, Visualization, Writing – Original Draft Preparation, Writing – Review & Editing

**Competing interests:** No competing interests were disclosed.

**Grant information:** I.I.B. is supported by the Canadian Institutes of Health Research, COVID-19 Rapid Research Funding Opportunity (02179-000). K.K. is the founder of BlueDot, a social enterprise that develops digital technologies for public health. K.K., A.W., A.T.B. and C.H. are employed at BlueDot. I.I.B. has consulted for BlueDot. T.d.O. and the NGS-SA is funded by the South African Medical Research Council (SAMRC), MRC SHIP and the Department of Science and Innovation (DSI) of South Africa. N.R.F. acknowledges support from a Wellcome Trust and Royal Society Sir Henry Dale Fellowship (204311/Z/16/Z) and a Medical Research Council-São Paulo Research Foundation CADDE partnership award (MR/S0195/1 and FAPESP 18/14389-0). V.H. was supported by the Biotechnology and Biological Sciences Research Council (BBSRC) [grant number BB/M010996/1]. M.U.G.K. acknowledges support from the Branco Weiss Fellowship and EU grant 874850 MOOD. The contents of this publication are the sole responsibility of the authors and do not necessarily reflect the views of the European Commission. O.G.P., J.P.M. and M.U.G.K. acknowledge support from the Oxford Martin School. AR acknowledges the support of the Wellcome Trust (Collaborators Award 206298/Z/17/Z – ARTIC network) and the European Research Council (grant agreement no. 725422 – ReservoirDOCS). A.OT is supported by the Wellcome Trust Hosts, Pathogens & Global Health Programme [grant number: grant.203783/Z/16/Z] and Fast Grants [award number: 2236]. COG-UK is supported by funding from the Medical Research Council (MRC) part of UK Research & Innovation (UKRI), the National Institute of Health Research (NIHR) and Genome Research Limited, operating as the Wellcome Sanger Institute. TFS acknowledges support from the Deutsche Forschungsgemeinschaft (SFB900, EXC2155 RESIST). SeqCOVID-SPAIN is supported by a grant from the Instituto de Salud Carlos III COV0020/00140.

**Copyright:** © 2021 O'Toole Á *et al.* This is an open access article distributed under the terms of the [Creative Commons Attribution License](#), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

**How to cite this article:** O'Toole Á, Hill V, Pybus OG *et al.* **Tracking the international spread of SARS-CoV-2 lineages B.1.1.7 and B.1.351/501Y-V2 with grinch [version 2; peer review: 3 approved]** Wellcome Open Research 2021, 6:121 <https://doi.org/10.12688/wellcomeopenres.16661.2>

**First published:** 19 May 2021, 6:121 <https://doi.org/10.12688/wellcomeopenres.16661.1>



**REVISED Amendments from Version 1**

We have updated the figures to amend some issues with proofing. We have added in some details of other excellent resources for SARS-CoV-2 international surveillance. Over on cov-lineages.org (which has had a facelift since time of publishing), we have also added in a resources page (<https://cov-lineages.org/resources.html>) that points the user to both internally developed and externally developed resources for SARS-CoV-2 lineage and variant tracking. Figure 1 and Figure 2 along with title were also updated.

**Any further responses from the reviewers can be found at the end of the article**

**Introduction**

In December 2020, routine genomic surveillance in the United Kingdom (UK)<sup>1</sup> reported a new and genetically distinct phylogenetic cluster of severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2) (variant VOC202012/01, lineage B.1.1.7). Preliminary analysis suggests that this lineage carries an unusually large number of genetic changes<sup>2</sup>. The earliest known cases of B.1.1.7 were sampled in southern England in late September 2020, and by December the lineage had spread to most UK regions and was growing rapidly<sup>3</sup>. In October 2020, a separate SARS-CoV-2 cluster (variant 501Y.V2, lineage B.1.351), which carried a different constellation of genetic changes, was detected by the Network for Genomic Surveillance in South Africa<sup>4,5</sup>. Both lineages carry mutations, especially in the virus spike protein, that may affect virus function, and both appear to have grown rapidly in relative frequency since their discovery. Early analyses of the spatial spread of SARS-CoV-2 highlights the potential for rapid virus dissemination through national and international travel<sup>6,7</sup>. Therefore continued genomic monitoring of lineages of concern is required.

To facilitate tracking of these lineages on an international scale, we developed a software tool *grinch* (global report investigating novel coronavirus haplotypes) that collates SARS-CoV-2 genomic data and epidemiological metadata. Resources such as *grinch* on cov-lineages.org can inform public health bodies and institutions around the world. Other excellent resources to track lineages and variants are available, including covariants.org, which tracks the spread of SARS-CoV-2 variants of interest, and outbreak.info, which gathers multiple sources of genetic and epidemiological data to track lineages. We include a non-exhaustive list of resources for tracking SARS-CoV-2 at <https://cov-lineages.org/resources.html>.

**Methods**

To better characterise the international distribution of lineages B.1.1.7 and B.1.351 we sourced SARS-CoV-2 sequences

from GISAID<sup>8,9</sup> and assigned lineages using pangolin (v2.1.6, <https://github.com/cov-lineages/pangolin>), which implements the nomenclature scheme described in Rambaut *et al.*,<sup>10</sup>. Genomes are assigned lineage B.1.1.7 if they exhibit at least 5 of the 17 mutations inferred to have arisen on the phylogenetic branch immediately ancestral to the cluster (Table 1)<sup>2</sup>; or to B.1.351 if they exhibit at least 5 of 9 lineage-associated mutations (Table 1)<sup>5</sup>. Lineage count and frequency data have been calculated daily using *grinch*. Using International Air Transport Association (IATA) travel data from October 2020, available through bluedot.global, we aggregated and collated the passenger volumes from international airports in London and South Africa to international destinations on same booking. Destinations with more than 5,000 passengers from London and more than 300 passengers from South Africa during the month of October are displayed on the cov-lineages.org website and in the underlying data for this publication<sup>11</sup>. *grinch*, with custom python modules that make use of geopandas v0.9, matplotlib v3.2 and seaborn v0.10, combines this information and produces reports with descriptive tables and figures that can be found at [https://cov-lineages.org/global\\_report.html](https://cov-lineages.org/global_report.html).

**Implementation**

All of the code underlying this daily lineage tracking web-report can be found at GitHub and Zenodo<sup>12</sup>. *grinch* is a python-based tool, the analysis pipeline of which is built on a snakemake backbone<sup>13</sup>. Every 24 hours a scheduled cron<sup>14</sup> task runs on our local servers. We download the latest data from GISAID and deduplicate based on sequence names. The sequences are assigned their most likely lineage using pangolin's latest version and model files. All processed metadata is available and maintained on the cov-lineages.org GitHub repository. To run *grinch*, the user must have access to a GISAID direct download key and a password and provide these within a configuration file for use. The command used to run *grinch* is `grinch -i grinch_config.yaml`, using the config file provided at doi:10.5281/zenodo.4640379<sup>15</sup>.

**Operation**

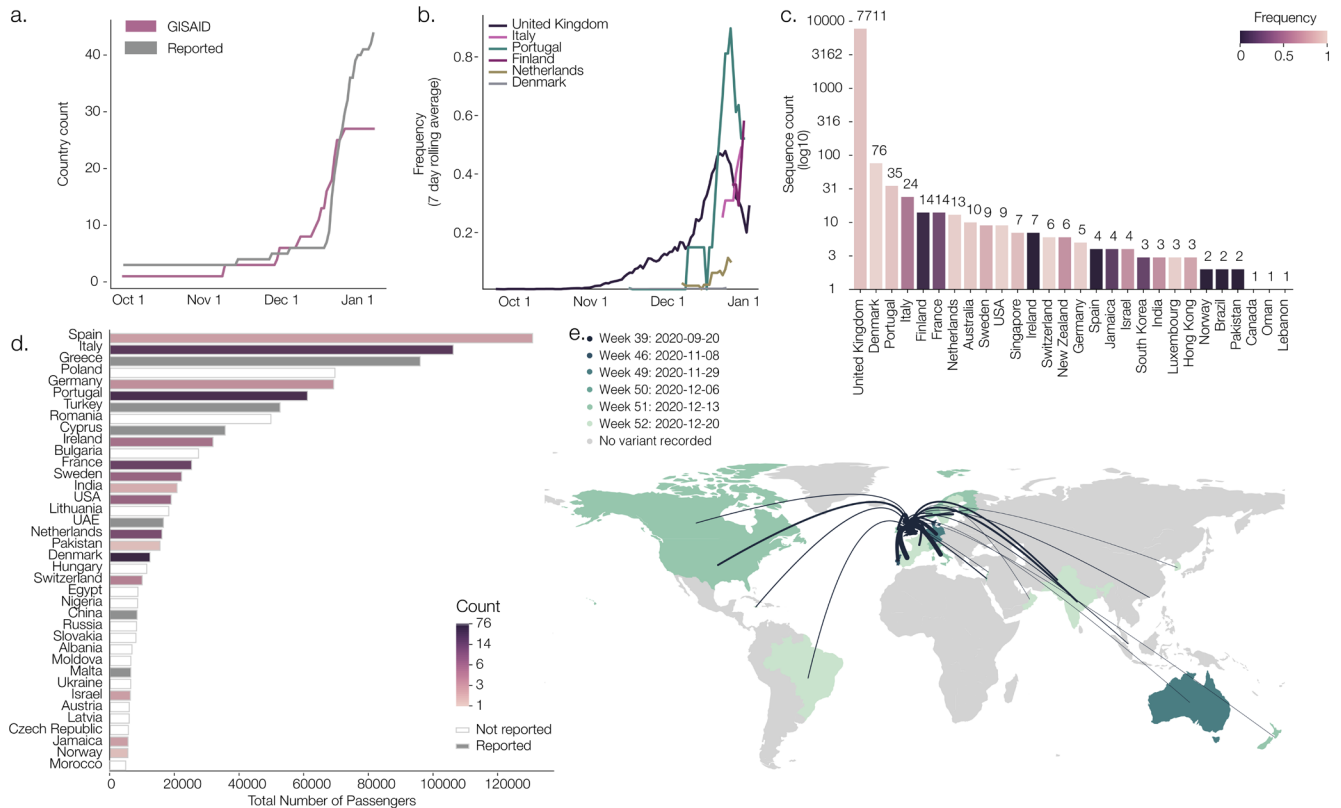
Most users will not run *grinch* themselves, instead all information and useful descriptive figures are provided daily on the web report. Users can navigate to cov-lineages.org in a web browser of choice to view the latest daily report.

**Results and discussion**

As of 7th Jan 2021, 45 countries had reported the presence of B.1.1.7 and 13 countries had reported B.1.351/501Y.V2. B.1.1.7 and B.1.351 genome sequences were available for 28 and 8 countries, respectively (Figure 1a, b, c)<sup>11</sup>. Although some countries report increases in the relative frequency of B.1.1.7, genome sequencing efforts vary considerably. Potential targeting of sequencing towards travelers from the UK could bias

**Table 1. Defining mutations for lineages of interest.**

Lineage	Defining mutations
B.1.1.7	orf1ab:T1001I; orf1ab:A1708D; orf1ab:I2230T; del:11288:9; del:21765:6; del:21991:3; S:N501Y; S:A570D; S:P681H; S:T716I; S:S982A; S:D1118H; Orf8:Q27*; Orf8:R52I; Orf8:Y73C; N:D3L; N:S235F
B.1.351/501Y-V2	E:P71L; N:T205I; orf1a:K1655N; S:D80A; S:D215G; S:K417N; S:E484K; S:N501Y; S:E484K



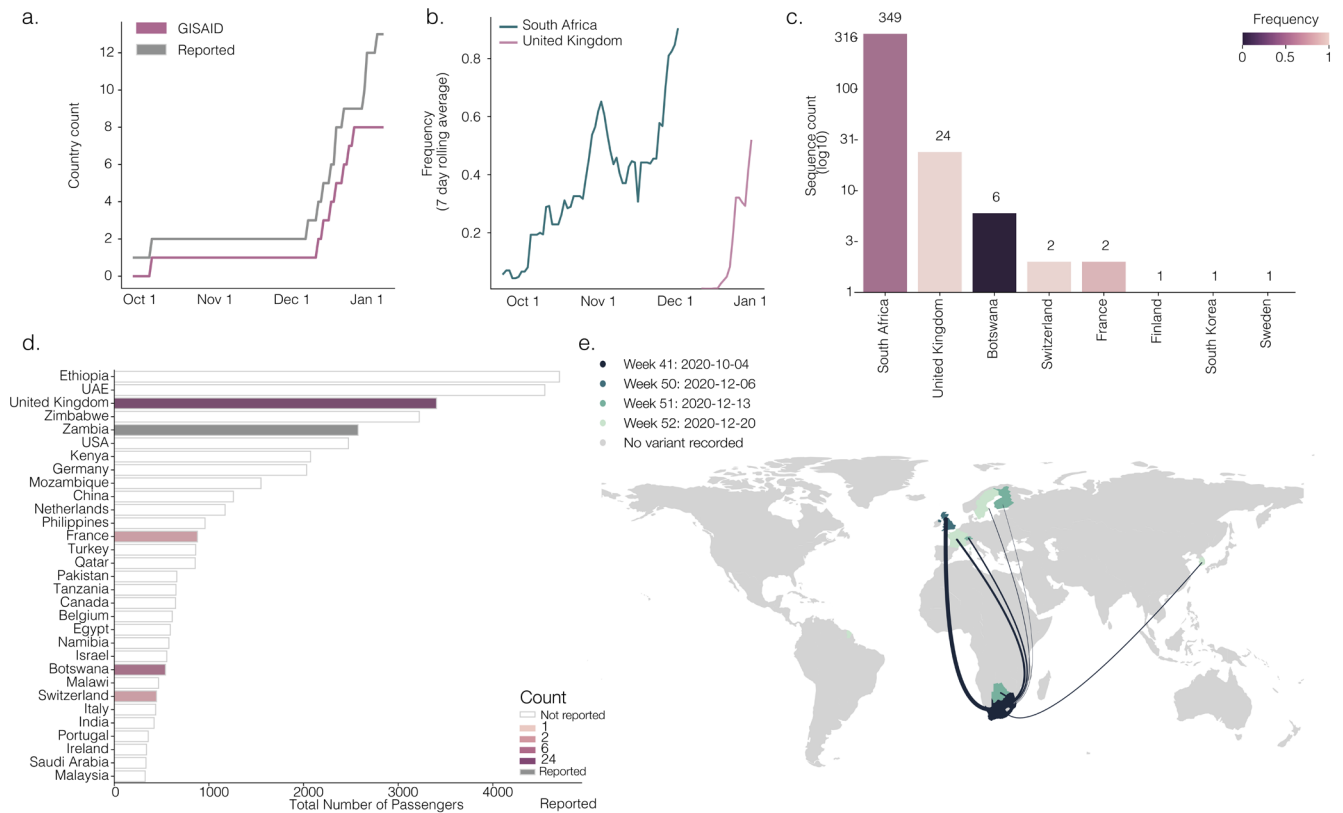
**Figure 1.** **a**) The cumulative number of countries with reports of lineage B.1.1.7 (grey line) and cumulative number of genomes of B.1.1.7 deposited in GISAID. **b**) Rolling seven-day average of the proportion of B.1.1.7 genomes in countries with more than ten sequences of the variant, and with more than ten days between the first B.1.1.7 sequence and the most recent one compared to all sampled genomes in that country. **c**) Number of sequences (log<sub>10</sub>) per country. Colour indicates the proportion of sequences that are classified as lineage B.1.1.7. **d**) Number of air travellers from major international London airports (Heathrow, Gatwick, Luton, City, Stansted, Southend) during October 2020. Colour indicates the number of sampled genomes of lineage B.1.1.7. Reported refers to countries that we found media reports stating there had been sequences of that particular lineage, but for which there were no sequences on GISAID. This is distinct from 'not reported' where there were no records found of that lineage in a given country. **e**) Map of international flights from major international London airports to countries with B.1.1.7 sequences. Colours indicate the date of earliest detection of B.1.1.7. in each country. The width of the lines indicates the number of flights. International Air Transport Association data used here account for ~90% of passenger travel itineraries on commercial flights, excluding transportation via unscheduled charter flights (the remainder is modelled using market intelligence). Data shown represents origin-destination journeys during October 2020. Routes to countries that have not yet detected B.1.1.7 and deposited data on GISAID are not included.

frequency estimates upwards (Figure 1b, c) and differing genome sharing policies and delays may also skew reporting estimates. The time between the initial collection date of a new variant sample in a country and the first availability of a corresponding virus genome on GISAID was, on average, 12 days (range 1–71).

The number of B.1.1.7 and B.1.351/501Y.V2 genome sequences reported in each country is a consequence of (i) the intensity of local genomic surveillance; (ii) the level of concern about new variant introductions; (iii) the volume of international travel among affected countries, and (iv) the amount of local transmission following the introduction of lineage from elsewhere. To explore these factors, we analysed the most recent available IATA travel data (October 2020). We collated the total number of origin-to-destination air journeys between major London

international airports and each country. The calculation was repeated for journeys originating in all international South African airports. We focussed on London and South Africa as they are the locations with the first reports and highest reported prevalence of lineages B.1.1.7 and B.1.351 respectively<sup>2,5</sup>. However, due to low SARS-CoV-2 genomic surveillance in many locations, we cannot reject the hypotheses that these lineages initially originated elsewhere. Figure 1d shows destinations receiving >5,000 travellers in October 2020 from the UK (Figure 2 shows destinations receiving >300 travellers from South Africa).

Of the countries that receive >5,000 travellers from London, 16 have sequenced B.1.1.7. Of the 45 countries that have identified B.1.1.7 (32 in travellers and 13 with local onward transmission), only 6 perform real-time routine genomic



**Figure 2.** **a)** Shows the cumulative number of countries with reports of lineage B.1.351 (black line) and cumulative number of genomes of B.1.351 deposited in GISAID. **b)** Rolling seven-day average of the proportion of B.1.351 genomes in countries with more than ten sequences of the variant, and with more than ten days between the first B.1.351 sequence and the most recent one compared to all sampled genomes in that country. **c)** Number of sequences (log<sub>10</sub>) per country. Colour indicates the proportion of sequences that are classified as lineage B.1.351. **d)** Number of air travellers from South Africa during October 2020. Colour indicates the number of sampled genomes of lineage B.1.351. Not reported refers to a given country having no record of B.1.351, and reported refers to countries that we found media reports but that country had no SARS-CoV-2 genomes shared on GISAID at that time. **e)** Map of international flights to countries with B.1.351 sequences. Colours indicate the date of earliest detection of B.1.351 in each country. The width of the lines indicates the number of flights. International Air Transport Association data used here account for ~90% of passenger travel itineraries on commercial flights, excluding transportation via unscheduled charter flights (the remainder is modelled using market intelligence). Data shown represents origin-destination journeys during October 2020. Routes to countries that have not yet detected B.1.351 and deposited data on GISAID are not included. >300 travellers from South Africa).

surveillance (Denmark, UK, Iceland, The Netherlands, Australia, Sweden), 3 have prioritised sequencing based on S-gene target failure tests<sup>16</sup>, 30 primarily targeted sequencing towards arriving travellers from the UK, and there was no information available for 10 (details at [https://github.com/cov-lineages/lineages-website/blob/master/\\_data/](https://github.com/cov-lineages/lineages-website/blob/master/_data/)). Of the 13 countries that have identified B.1.351 (four with local onward transmission including South Africa), 4 perform routine sequencing (South Africa, UK, Botswana, Australia), 6 target sequencing of travellers, and there was no information available for 3. Consequently, there is no clear relationship between number of sequences reported and flight numbers, but rather reflects the current genomic surveillance effort. For example, in September, the UK sequenced ~13% of its reported cases and Denmark sequenced ~21%. In comparison, Israel sequenced ~0.002% of its cases during the same period<sup>17,18</sup>.

Our study has several limitations. The passenger flight data do not include recent changes to holiday travel, and recent restrictions on travel from the UK and South Africa is not reflected in the mobility data. Further, flight data may not accurately reflect the final destination if multiple tickets are purchased.

The discovery and rapid spread of B.1.1.7 and B.1.351/501Y.V2 highlights the importance of real-time and open data for tracking the spread of SARS-CoV-2 and for informing future public health interventions and travel advice.

## Data availability

### Underlying data

Zenodo: Accession IDs included in publication Tracking the international spread of SARS-CoV-2 lineages B.1.1.7 and B.1.351/501Y.V2. <https://doi.org/10.5281/zenodo.4642401>.



This project contains the following underlying data:

- Accession IDs of B.1.1.7 and B.1.351 genome sequences included in report up until January 7<sup>th</sup>, 2021. All accession IDs link to data on the GISAID repository, <http://doi.org/10.17616/R3Q59F>. These data are available under the terms of the [GISAID EpiFlu™ Database Access Agreement](#).

Zenodo: [cov-lineages.org](https://doi.org/10.5281/zenodo.4640140) website. <https://doi.org/10.5281/zenodo.4640140><sup>11</sup>.

This project contains the following underlying data:

- Website data archived at time of publication

Data are available under the terms of the [Creative Commons Attribution 4.0 International license](#) (CC-BY 4.0).

#### Extended data

Zenodo: Supplementary materials with group affiliations for Tracking the international spread of SARS-CoV-2 lineages B.1.1.7 and B.1.351/501Y-V2. <https://doi.org/10.5281/zenodo.4704471><sup>19</sup>.

This project contains the following extended data:

- Supplementary materials with group authorship affiliations and full acknowledgements.

Data are available under the terms of the [Creative Commons Attribution 4.0 International license](#) (CC-BY 4.0).

#### Software availability

- **Software available from:** [https://cov-lineages.org/global\\_report.html](https://cov-lineages.org/global_report.html)
- **Source code available from:** <https://github.com/cov-lineages/grinch>

- **Archived source code at time of publication:** <https://doi.org/10.5281/zenodo.4640037><sup>12</sup>; <https://doi.org/10.5281/zenodo.4640379><sup>15</sup>

- **Licenses:** GNU General Public License v3.0; Creative Commons Attribution 4.0 International license (CC-BY 4.0).

#### Author contributions

A.OT., O.G.P., J.P.M., N.R.F., A.R., M.U.G.K. conceived the study. A.OT, V.H., M.U.G.K. O.G.P., A.R. wrote the first draft. A.OT, V.H., M.U.G.K., A.R., AW, conducted data analysis. S.T., T.Salloum, G.M., J.K., J.G., J.d.L., X.R., M.S, N.Freed, C.P., P.P., A.D., R.V., T.F.S., L.S., T.Stadler, A.P., A.B., D.G.d.V., S.B-S., V.B., J.I., S.D., J.P.G., N.Z., M.M., O.M., T.Seemann, N.S., B.H., M.Sait, A.A., J.D., M.G., W.R., I.D., S.S., J.M., L.L., C.K.L., M.C., D.L., A.B., L.J., K.S.G., E.L-N., R.S., B.M., M.Koopmans, M.B., V.S.R., S.P., T.Smura, A.H., S.K., M.U., M.Salman, B.B., M.R., C.D., T.W., O.S., D.E., C.R., H.V., A.P. contributed to the genomic dataset and facilitated data and sample availability. All authors interpreted the data and contributed to writing.

#### Acknowledgements

An earlier version of this article can be found on Virological (url: <https://virological.org/t/tracking-the-international-spread-of-sars-cov-2-lineages-b-1-1-7-and-b-1-351-501y-v2/592>).

We thank Norelle Sherry, Benjamin Howden and Michelle Sait for their contribution to sequencing in Australia. We also include full acknowledgements and details of group authorships at <https://doi.org/10.5281/zenodo.4704471><sup>19</sup>. We would also like to extend our gratitude to everyone involved in the global sequencing effort.

#### References

1. COVID-19 Genomics UK (COG-UK) consortium [contact@cogconsortium.uk](mailto:contact@cogconsortium.uk): **An integrated national scale SARS-CoV-2 genomic surveillance network**. *Lancet Microbe*. 2020; 1(3): e99–100. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
2. Rambaut A, Loman N, Pybus O, et al.: **Preliminary genomic characterisation of an emergent SARS-CoV-2 lineage in the UK defined by a novel set of spike mutations**. 2020; published online Dec 18. (accessed Jan 8, 2021). [Reference Source](#)
3. Volz E, Mishra S, Chand M, et al.: **Transmission of SARS-CoV-2 Lineage B.1.1.7 in England: Insights from linking epidemiological and genetic data**. *bioRxiv*. 2021. [Publisher Full Text](#)
4. Msomi N, Mlisana K, de Oliveira T: **A genomics network established to respond rapidly to public health threats in South Africa**. *Lancet Microbe*. 2020; 1(6): e229–30. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
5. Tegally H, Wilkinson E, Giovanetti M, et al.: **Emergence and rapid spread of a new severe acute respiratory syndrome-related coronavirus 2 (SARS-CoV-2) lineage with multiple spike mutations in South Africa**. *bioRxiv*. 2020. [Publisher Full Text](#)
6. du Plessis L, McCrone JT, Zarebski AE, et al.: **Establishment and lineage dynamics of the SARS-CoV-2 epidemic in the UK**. *Science*. 2021; 371(6530): 708–712. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
7. Lu J, du Plessis L, Liu Z, et al.: **Genomic Epidemiology of SARS-CoV-2 in Guangdong Province, China**. *Cell*. 2020; 181(5): 997–1003.e9. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
8. Elbe S, Buckland-Merrett G: **Data, disease and diplomacy: GISAID's innovative contribution to global health**. *Glob Chall*. 2017; 1(1): 33–46. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
9. O'Toole A: **Accession IDs included in publication Tracking the international spread of SARS-CoV-2 lineages B.1.1.7 and B.1.351/501Y-V2 [Data set]**. *Zenodo*. 2021. <http://www.doi.org/10.5281/zenodo.4642401>
10. Rambaut A, Holmes EC, O'Toole Á, et al.: **A dynamic nomenclature proposal for SARS-CoV-2 lineages to assist genomic epidemiology**. *Nat Microbiol*. 2020; 5(11): 1403–7. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
11. O'Toole A: **cov-lineages.org website**. *Zenodo*. 2021. <http://www.doi.org/10.5281/zenodo.4640140>
12. O Toole A, Hill V: **grinch**. *Zenodo*. 2021. <http://www.doi.org/10.5281/zenodo.4640037>
13. Mölder F, Jablonski KP, Letcher B, et al.: **Sustainable data analysis with Snakemake [version 1; peer review: 1 approved, 1 approved with reservations]**. *F1000Res*. 2021; 10: 33. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
14. Reznick L: **Using cron and crontab**. *Sys Admin*. 1993; 2(4): 29–32.
15. O'Toole A: **grinch\_config.yaml [Data set]**. *Zenodo*. 2021. <http://www.doi.org/10.5281/zenodo.4640379>

16. Bal A, Destras G, Gaymard A, *et al.*: **Two-step strategy for the identification of SARS-CoV-2 variants co-occurring with spike deletion H69-V70, Lyon, France, August to December 2020.** *bioRxiv.* 2020.  
[PubMed Abstract](#) | [Publisher Full Text](#)
17. Hasell J, Mathieu E, Beltekian D, *et al.*: **A cross-country database of COVID-19 testing.** *Sci Data.* 2020; 7(1): 345.  
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
18. Dong E, Du H, Gardner L: **An interactive web-based dashboard to track COVID-19 in real time.** *Lancet Infect Dis.* 2020; 20(5): 533–4.  
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
19. O'Toole A: **Supplementary materials with group affiliations for Tracking the international spread of SARS-CoV-2 lineages B.1.1.7 and B.1.351/501Y-V2.** *Zenodo.* 2021.  
<http://www.doi.org/10.5281/zenodo.4704471>

# Open Peer Review

Current Peer Review Status:   

---

## Version 1

Reviewer Report 10 June 2021

<https://doi.org/10.21956/wellcomeopenres.18372.r43967>

© 2021 Githinji G. This is an open access peer review report distributed under the terms of the [Creative Commons Attribution License](#), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.



**George Githinji** 

Epidemiology and Demography Department, Kenya Medical Research Institute (KEMRI) - Wellcome Trust Research Programme, Kilifi, Kenya

The article by O'Toole 2021 *et al.* describes a bioinformatics tool for the analysis of SARS-CoV-2 sequence data. The article is concise, and the relevant details have been considered. For example, the software and source code is available and well documented. The tool has shown great utility in public health based on its application in tracking and describing two SARS-CoV-2 variants of global concern.

Some minor comments below:

1. The transition from an article of public health importance to a software tool is abrupt. I think a paragraph or a link aimed at orientating the audience would be useful.
2. It would be useful to outline the special niche that the tool occupies or the gaps it fills relative to similar utilities and webpages such as [covariants.org](https://covariants.org) and [outbreak.info](https://outbreak.info).
3. The readme file at <https://github.com/cov-lineages/grinch> lacks full installation documentation. An introductory paragraph of the tool and its utility would also be useful. The scripts directory could be better organised by separating the snakemake files from the regular python files. I would image a workflows dir and scripts dir

**Is the rationale for developing the new software tool clearly explained?**

Yes

**Is the description of the software tool technically sound?**

Yes

**Are sufficient details of the code, methods and analysis (if applicable) provided to allow replication of the software development and its use by others?**

Yes

**Is sufficient information provided to allow interpretation of the expected output datasets and any results generated using the tool?**

Yes

**Are the conclusions about the tool and its performance adequately supported by the findings presented in the article?**

Yes

**Competing Interests:** No competing interests were disclosed.

**Reviewer Expertise:** bioinformatics, molecular epidemiology

**I confirm that I have read this submission and believe that I have an appropriate level of expertise to confirm that it is of an acceptable scientific standard.**

Author Response 10 Sep 2021

**Áine Niamh O'Toole**, Institute of Evolutionary Biology, University of Edinburgh, Edinburgh, UK

Thank you for your review and time, and apologies for taking so long to respond!

**The transition from an article of public health importance to a software tool is abrupt. I think a paragraph or a link aimed at orientating the audience would be useful.**

We have added in a paragraph that bridges the public health information and the resource information.

**It would be useful to outline the special niche that the tool occupies or the gaps it fills relative to similar utilities and webpages such as covariants.org and outbreak.info.**

The linking paragraph also discusses other resources such as outbreak.info and covariants.org- we now also provide a non-exhaustive list of a several other resources on the cov-lineages.org website

**The readme file at <https://github.com/cov-lineages/grinch> lacks full installation documentation. An introductory paragraph of the tool and its utility would also be useful. The scripts directory could be better organised by separating the snakemake files from the regular python files. I would image a workflows dir and scripts dir**

We have updated the readme on the grinch repository and updated usage in that it can be run in full analysis pipeline mode or in report only mode. We have also added in a brief description of the grinch setup and where the macro data is hosted.

**Competing Interests:** No competing interests were disclosed.

Reviewer Report 03 June 2021

<https://doi.org/10.21956/wellcomeopenres.18372.r43964>

© 2021 F. Brito A. This is an open access peer review report distributed under the terms of the [Creative Commons Attribution License](#), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.



**Anderson F. Brito** 

Department of Epidemiology of Microbial Diseases, Yale School of Public Health, New Haven, CT, USA

Thank you for developing these tools for daily tracking of SARS-CoV-2 lineages.

**Some comments about the tool and its functionalities:**

- In the manuscript, the authors mention that users with access to GISAID direct download could run this pipeline locally, and generate their own reports. Can this tool be adapted to display genomic results for tracking national spread, or even state-level spread of lineages?
- If this pipeline is intended to be constantly executed locally by users, it would be helpful to provide more information about how to install and run the pipeline, including reference to example input and output files. I have tried to run the pipeline using my GISAID data provision credentials, but that was not successful, as I ran into errors for which I could not find a solution online (GitHub and Zenodo).
- About the online reports, increasing the font size in the plots being displayed (bar, curves, etc) would make labels and legends more intelligible, and improving the readability of their content.
- About the flight data, why only flight counts from October are shown? Are these data only used for tracking the potential spread in early stages of viral emergence, or do you see other uses for such data?

**Concerning the manuscript, a few minor points:**

- The colour gradient in the legend of Figure 1 is incomplete and does not go from 1 to 76. I think it must be just a formatting issue.
- How was the "reported" cases shown in Figures 1 and 2 detected? By differential PCR? I know that applies to B.1.1.7, but what about B.1.351?
- The legend in Figure 2 refers to "B.1.1.7" sequences, while the figure shows "B.1.351" sequences. It must be a typo.

**Is the rationale for developing the new software tool clearly explained?**

Yes

**Is the description of the software tool technically sound?**

Partly

**Are sufficient details of the code, methods and analysis (if applicable) provided to allow**



**replication of the software development and its use by others?**

Partly

**Is sufficient information provided to allow interpretation of the expected output datasets and any results generated using the tool?**

Partly

**Are the conclusions about the tool and its performance adequately supported by the findings presented in the article?**

Yes

**Competing Interests:** No competing interests were disclosed.

**Reviewer Expertise:** Virology, Bioinformatics, Evolution, Epidemiology.

**I confirm that I have read this submission and believe that I have an appropriate level of expertise to confirm that it is of an acceptable scientific standard.**

Author Response 10 Sep 2021

**Áine Niamh O'Toole**, Institute of Evolutionary Biology, University of Edinburgh, Edinburgh, UK

- **In the manuscript, the authors mention that users with access to GISAID direct download could run this pipeline locally, and generate their own reports. Can this tool be adapted to display genomic results for tracking national spread, or even state-level spread of lineages?**

We've added in an option to grinch with the --analysis flag, that can either run the whole GISAID-processing pipeline, or just generate a set of reports from a metadata file. It currently would need resolution to remain at the country level, but we like the idea of doing a more local report and will work towards implementing something like this in the future.

- **If this pipeline is intended to be constantly executed locally by users, it would be helpful to provide more information about how to install and run the pipeline, including reference to example input and output files. I have tried to run the pipeline using my GISAID data provision credentials, but that was not successful, as I ran into errors for which I could not find a solution online (GitHub and Zenodo).**

We have supplied some more information on how to install on the GitHub repository and run the pipeline, however the tool isn't necessarily intended for users to run locally themselves as we process the data on this end and share the macro count data on GitHub.

- **About the online reports, increasing the font size in the plots being displayed (bar, curves, etc) would make labels and legends more intelligible, and improving the readability of their content.**

We have made the axes longer to account for more countries, however are working towards re-implementing these reports in javascript so they can be interactive and more responsive to the browser size.

- **About the flight data, why only flight counts from October are shown? Are these data only used for tracking the potential spread in early stages of viral emergence, or do you see other uses for such data?**

October related to the date of early spread of both lineages described in the text and was due to limitations of access to data. We hope to continue to develop this resource and supply more recent dates that track over time.

**Concerning the manuscript, a few minor points:**

- **The colour gradient in the legend of Figure 1 is incomplete and does not go from 1 to 76. I think it must be just a formatting issue.**

This was a proofs issue and we've hopefully rectified this now.

- **How was the "reported" cases shown in Figures 1 and 2 detected? By differential PCR? I know that applies to B.1.1.7, but what about B.1.351?**

This was the set of manually curated media reports that we were tracking at the time.

- **The legend in Figure 2 refers to "B.1.1.7" sequences, while the figure shows "B.1.351" sequences. It must be a typo.**

Thank you we have now fixed this in the figure.

**Competing Interests:** No competing interests were disclosed.

Author Response 10 Sep 2021

**Áine Niamh O'Toole**, Institute of Evolutionary Biology, University of Edinburgh, Edinburgh, UK

Also- thank you for your time reviewing this manuscript (we do really appreciate it!) and apologies for taking so long to respond.

**Competing Interests:** No competing interests were disclosed.

Reviewer Report 27 May 2021

<https://doi.org/10.21956/wellcomeopenres.18372.r43966>

© 2021 Lanfear R. This is an open access peer review report distributed under the terms of the [Creative Commons Attribution License](#), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.



**Rob Lanfear** 

Department of Ecology and Evolution, Research School of Biology, Australian National University, Canberra, ACT, Australia

This article describes a software tool, *grinch*, that can be used to produce automated reports on SARS-CoV-2 lineages. The authors apply it to two lineages of concern in the article, and also highlight that the main utility of *grinch* is not in static one-off reports, but in regularly updated reports available at [https://cov-lineages.org/global\\_report.html](https://cov-lineages.org/global_report.html).

The paper clearly describes the software and demonstrates its utility. I'd like to commend the authors for putting this tool and the associated website together so quickly, for maintaining both to a very high standard, for making sure that all of the work is open and reproducible, and for the huge amount of work and enormous collaborative effort that has gone into this clear and concise report.

I have no serious reservations about the software tool or the data, analyses, or conclusions presented in the manuscript. The software is clear, open-source, sufficiently documented, and almost all of the proposed utility is presented on a clear and regularly updated website. The manuscript is clearly written, well researched, concise, and the conclusions are well justified by the analyses.

Of course, I do have a few comments, some of which I hope might be useful in improving the paper and/or the website.

Minor comments on the manuscript:

1. I felt there was some tension in this article about whether it's a software note or a public health report. The title suggests the latter, but much of the article (and the article type of "Software Tool Article") suggests the former. Most of this tension for me as a reader came from looking at the title, which has no mention of software, so I think sets up expectations that differ from what is then provided (quite reasonably) in the paper. A very simple way to address this would be to start the title with "Using *grinch* to track..." or to end it with "... using *grinch*".
2. Similar to point 1, the abstract doesn't actually mention '*grinch*' or [https://cov-lineages.org/global\\_report.html](https://cov-lineages.org/global_report.html). It would seem clearer to me to incorporate in the abstract the framing that this article presents a generally applicable software tool, demonstrated on two lineages of concern.
3. I would like to see some mention of related efforts somewhere in the report. A full detailed comparison is neither warranted nor useful here because all such websites can and should change regularly, but a couple of sentences comparing [cov-lineages.org](https://cov-lineages.org) to sites like [outbreak.info](https://outbreak.info) and [covariants.org](https://covariants.org) would be very useful. At a minimum, it seems useful to list the similar sites the authors are aware of, if only because the fact one can see similar patterns presented on those sites serves as a useful validation of the software presented in this paper.
4. Given the situation, this is a desirable, not a requirement, but I'd love to see some unit tests on the GitHub repo. It seems potentially important to have this when the intention is to

produce daily updates for public health. (Though I note that getting the same end result from completely independent implementations on other sites is probably worth more than a lot of unit tests).

5. I struggled with Figure 1D. It wasn't clear to me what 'reported' and 'not reported' mean. And the legend makes it really hard to figure out how colours map to counts.
6. It's stated that there is no correlation between the numbers of sequences and flight numbers. It would be nice to see the scatter plot for this (maybe as an inset to figure 1D?), as well as the effect size and p-value of a suitable model.
7. Following from the previous point, the explanation for the lack of a correlation with absolute numbers seems reasonable. But it still seems to me that flight numbers could correlate with the frequency of B117 at a fixed time interval from the first detected case in a given locality (thus somewhat factoring out sequencing effort in the locality). Is it possible to add this analysis?
8. Please add installation instructions to the GitHub repo

Minor comments on [https://cov-lineages.org/global\\_report.html](https://cov-lineages.org/global_report.html):

1. Figure 3 for each lineage is a map of sequence counts by region. I find the legend here completely baffling. All it states is grey=No variant (that makes sense), pink = 1 sequence (that makes sense too), and purple = 'Max sequences'. I have no idea what to make of this. How many is 'Max', and how am I supposed to quantitatively interpret intermediate colours to pink and purple? It's so obvious I'm certain there are good reasons why this isn't already done, but it does seem like a continuous colour scale is what *should* be used here. Similar to the scale in Figure 2 (grey for no data, shades of green nicely spaced and annotated for different values of a continuous variable).
2. For the widespread lineages like B.1.1.7, there's a lot of overplotting on Figures 4 and 5, which make the counts and the country names very difficult to read. This could be addressed by just making the figures larger.
3. The table of links to news reports is absolutely wonderful. Would it be possible to include a button here to allow users to suggest additional news links? (I assume there's an existing mechanism for doing this, but I couldn't find one, so if not maybe just a link to a github issue with (potentially) a pre-filled title and required information would help?)

Really minor comments about the manuscript:

1. The first use of IATA (first para of the methods) is missing "International", i.e. it says "Using Air Transport...".
2. The second use of IATA (second para of results) does not need to be spelled out.
3. Figure 1A seems like it is missing a second Y axis for the number of GISAID genomes reported.
4. In the PDF version and the HTML version it seems that new lines were added wherever

there was a '>', e.g. '>5,000 travellers in October' and '>300 travellers from South Africa' both (erroneously?) start on new lines.

**Is the rationale for developing the new software tool clearly explained?**

Yes

**Is the description of the software tool technically sound?**

Yes

**Are sufficient details of the code, methods and analysis (if applicable) provided to allow replication of the software development and its use by others?**

Yes

**Is sufficient information provided to allow interpretation of the expected output datasets and any results generated using the tool?**

Yes

**Are the conclusions about the tool and its performance adequately supported by the findings presented in the article?**

Yes

**Competing Interests:** I am a paid consultant to GISAID, the database on which much of the data analysed in this article is hosted.

**Reviewer Expertise:** Phylogenetics, molecular evolution, bioinformatics. I have a passing familiarity with SARS-CoV-2 data analysis.

**I confirm that I have read this submission and believe that I have an appropriate level of expertise to confirm that it is of an acceptable scientific standard.**

Author Response 10 Sep 2021

**Áine Niamh O'Toole**, Institute of Evolutionary Biology, University of Edinburgh, Edinburgh, UK

Thank you so much for your detailed review and our apologies for not responding sooner. All your comments were well justified and fair. We've responded to and believe have addressed your comments and concerns below.

*I felt there was some tension in this article about whether it's a software note or a public health report. The title suggests the latter, but much of the article (and the article type of "Software Tool Article") suggests the former. Most of this tension for me as a reader came from looking at the title, which has no mention of software, so I think sets up expectations that differ from what is then provided (quite reasonably) in the paper. A very simple way to address this would be to start the title with "Using grinch to track..." or to end it with "... using grinch".*

I completely agree, it's a very fair point! We had originally submitted the article as a public



health report but this didn't fit with the Wellcome Open Research journal remit, so resubmitted under Software Tool. I have changed the title as suggested to with with "using grinch".

*Similar to point 1, the abstract doesn't actually mention 'grinch' or [https://cov-lineages.org/global\\_report.html](https://cov-lineages.org/global_report.html). It would seem clearer to me to incorporate in the abstract the framing that this article presents a generally applicable software tool, demonstrated on two lineages of concern.*

Our abstract and introduction now both contain reference to grinch and the reports at cov-lineages.org

*I would like to see some mention of related efforts somewhere in the report. A full detailed comparison is neither warranted nor useful here because all such websites can and should change regularly, but a couple of sentences comparing cov-lineages.org to sites like outbreak.info and covariants.org would be very useful. At a minimum, it seems useful to list the similar sites the authors are aware of, if only because the fact one can see similar patterns presented on those sites serves as a useful validation of the software presented in this paper.*

We have added in a short paragraph about these resources and a link to a more extensive (but definitely non-exhaustive) list of resources.

*Given the situation, this is a desirable, not a requirement, but I'd love to see some unit tests on the GitHub repo. It seems potentially important to have this when the intention is to produce daily updates for public health. (Though I note that getting the same end result from completely independent implementations on other sites is probably worth more than a lot of unit tests).*

Since publication, we've re-worked the back-end analysis pipeline and the GISAID data is now processed with the datapipe pipeline written by Rachel Colquhoun (<https://github.com/COG-UK/datapipe>). The reports and webpages are still generated with grinch, however the main data processing steps are now done with the robust datapipe pipeline.

*It's stated that there is no correlation between the numbers of sequences and flight numbers. It would be nice to see the scatter plot for this (maybe as an inset to figure 1D?), as well as the effect size and p-value of a suitable model. Following from the previous point, the explanation for the lack of a correlation with absolute numbers seems reasonable. But it still seems to me that flight numbers could correlate with the frequency of B117 at a fixed time interval from the first detected case in a given locality (thus somewhat factoring out sequencing effort in the locality). Is it possible to add this analysis?*

We have amended to state there is no clear relationship, rather than correlation.

*Please add installation instructions to the GitHub repo.*

We have added updated usage and install instructions, and a description of the behaviour to the repository, at <https://github.com/cov-lineages/grinch/blob/main/README.md>

**Minor comments on [https://cov-lineages.org/global\\_report.html](https://cov-lineages.org/global_report.html):**

1. **Figure 3 for each lineage is a map of sequence counts by region. I find the legend here completely baffling. All it states is grey=No variant (that makes sense), pink = 1 sequence (that makes sense too), and purple = 'Max sequences'. I have no idea what to make of this. How many is 'Max', and how am I supposed to quantitatively interpret intermediate colours to pink and purple? It's so obvious I'm certain there are good reasons why this isn't already done, but it does seem like a continuous colour scale is what should be used here. Similar to the scale in Figure 2 (grey for no data, shades of green nicely spaced and annotated for different values of a continuous variable).**

We have since amended the legend for the report.

2. **For the widespread lineages like B.1.1.7, there's a lot of overplotting on Figures 4 and 5, which make the counts and the country names very difficult to read. This could be addressed by just making the figures larger.**

We have made these figures larger, but recognise that these plots may be better displayed by having a top 20 country barchart (like in <https://cov-lineages.org/lineage.html?lineage=B.1.1.7>) and intend to adopt an interactive, more attractive global report in the future.

3. **The table of links to news reports is absolutely wonderful. Would it be possible to include a button here to allow users to suggest additional news links? (I assume there's an existing mechanism for doing this, but I couldn't find one, so if not maybe just a link to a github issue with (potentially) a pre-filled title and required information would help?)**

We really like this idea and will work towards implementing it!

**Competing Interests:** No competing interests were disclosed.