

# 1                    **Untangling introductions and persistence in COVID-19 resurgence in Europe**

2  
3  
4 Philippe Lemey<sup>1,2</sup>, Nick Ruktanonchai<sup>3,4</sup>, Samuel L. Hong<sup>1</sup>, Vittoria Colizza<sup>5</sup>, Chiara Poletto<sup>5</sup>, Frederik Van  
5 den Broeck<sup>1,6</sup>, Mandev S. Gill<sup>1</sup>, Xiang Ji<sup>7</sup>, Anthony Levasseur<sup>8</sup>, Bas B. Oude Munnink<sup>9</sup>, Marion  
6 Koopmans<sup>9</sup>, Adam Sadilek<sup>10</sup>, Shengjie Lai<sup>3</sup>, Andrew J. Tatem<sup>3</sup>, Guy Baele<sup>1</sup>, Marc A. Suchard<sup>11,12,13</sup>, Simon  
7 Dellicour<sup>1,14</sup>.

8  
9  
10 <sup>1</sup>Department of Microbiology, Immunology and Transplantation, Rega Institute, KU Leuven, Leuven, Belgium.

11 <sup>2</sup>Global Virus Network (GVN), Baltimore, MD, USA.

12 <sup>3</sup>WorldPop, School of Geography and Environmental Science, University of Southampton, Southampton SO17 1BJ,  
13 UK.

14 <sup>4</sup>Population Health Sciences, Virginia Tech, Blacksburg, VA, USA.

15 <sup>5</sup>INSERM, Sorbonne Université, Institut Pierre Louis d'Epidémiologie et de Santé Publique IPLESP, F75012 Paris,  
16 France.

17 <sup>6</sup>Department of Biomedical Sciences, Institute of Tropical Medicine, Antwerp, Belgium.

18 <sup>7</sup>Department of Mathematics, School of Science & Engineering, Tulane University, New Orleans, LA, USA

19 <sup>8</sup>UMR MEPHI (Microbes, Evolution, Phylogeny and Infections), Aix-Marseille Université (AMU) and Institut  
20 Universitaire de France (IUF), Marseille, France.

21 <sup>9</sup>Department of Viroscience, WHO Collaborating Centre for Arbovirus and Viral Hemorrhagic Fever Reference and  
22 Research, Erasmus MC, Rotterdam, The Netherlands.

23 <sup>10</sup>Google, Mountain View, CA, USA.

24 <sup>11</sup>Department of Biomathematics, David Geffen School of Medicine, University of California Los Angeles, Los Angeles,  
25 CA 90095, USA.

26 <sup>12</sup>Department of Biostatistics, Fielding School of Public Health, University of California Los Angeles, Los Angeles, CA  
27 90095, USA.

28 <sup>13</sup>Department of Human Genetics, David Geffen School of Medicine, University of California Los Angeles, Los Angeles,  
29 CA 90095, USA.

30 <sup>14</sup>Spatial Epidemiology Lab (SpELL), Université Libre de Bruxelles, CP160/12, 50 av. FD Roosevelt, 1050 Bruxelles,  
31 Belgium.

32 **Summary paragraph**

33 Following the first wave of SARS-CoV-2 infections in spring 2020, Europe experienced a resurgence of the  
34 virus starting late summer that was deadlier and more difficult to contain <sup>1</sup>. Relaxed intervention  
35 measures and summer travel have been implicated as drivers of the second wave <sup>2</sup>. Here, we build a  
36 phylogeographic model to evaluate how newly introduced lineages, as opposed to the rekindling of  
37 persistent lineages, contributed to the COVID-19 resurgence in Europe. We inform this model using  
38 genomic, mobility and epidemiological data from 10 European countries and estimate that in many  
39 countries over half of the lineages circulating in late summer resulted from new introductions since June  
40 15<sup>th</sup>. The success in onward transmission of newly introduced lineages was negatively associated with  
41 local COVID-19 incidence during this period. The pervasive spread of variants in summer 2020 highlights  
42 the threat of viral dissemination when restrictions are lifted, and this needs to be carefully considered by  
43 strategies to control the current spread of variants that are more transmissible and/or evade immunity.  
44 Our findings indicate that more effective and coordinated measures are required to contain spread  
45 through cross-border travel even as vaccination is reducing disease burden.

46

47 **Keywords:** COVID-19, SARS-CoV-2, Europe, second wave, phylogeography, international mobility

48 Upon successfully curbing transmission in spring 2020, many European countries witnessed a resurgence  
49 in COVID-19 cases in late summer. The number of COVID-19 infections increased rapidly, and by the end  
50 of October, it was clear that the continent was deep into a second epidemic wave. This forced  
51 governments to reimpose lockdowns and social restrictions in an effort to contain the resurgence. While  
52 these measures reduced infection rates across Europe<sup>3</sup>, several countries witnessed a stabilization at high  
53 levels or even a new surge in infections. The spread of more transmissible variants, in particular B.1.1.7  
54 (alpha variant or 20I/501Y.V1<sup>4</sup>), which was first identified in the United Kingdom (UK), has considerably  
55 exacerbated the challenge to contain COVID-19.

56

57 Already early on in the pandemic, modelling studies warned about new waves due to partial relaxation of  
58 restrictions<sup>5</sup> or seasonal variations<sup>6</sup>. By mid-April, the European Commission constructed a roadmap to  
59 lifting coronavirus containment measures<sup>7</sup>, recommending a cautious and coordinated manner to revive  
60 social and economic activities. However, the early start of the devastating second wave demonstrated  
61 that there was insufficient adherence to these measured recommendations. Cross-border travel, and  
62 mass tourism in particular, has been implicated as a major instigator of the second wave. Genomic  
63 surveillance demonstrated that a new variant (lineage B.1.177<sup>8</sup>, 20E (EU1) [nextstrain.org]), which  
64 emerged in Spain in early summer, has spread to multiple locations in Europe<sup>2</sup>. While this variant quickly  
65 grew into the dominant circulating SARS-CoV-2 strain in several countries, it did not appear to be  
66 associated with a higher intrinsic transmissibility<sup>2</sup>.

67

68 Although it appears clear that travel considerably contributed to the second wave in Europe, it remains  
69 challenging to assess how it may have restructured and reignited the epidemic in the different European  
70 countries. Even without resuming travel, relaxing containment measures when low-level transmission is  
71 ongoing risks the proliferation of locally circulating strains. Phylodynamic analyses may provide insights  
72 into the relative importance of persistence versus the introduction of new lineages, but such analyses are  
73 complicated for SARS-CoV-2 for different reasons. Phylogenetic reconstructions may be poorly resolved  
74 due to the relatively limited SARS-CoV-2 sequence diversity<sup>9</sup>. This is further confounded by the degree of  
75 genetic mixing that can be expected from unrestricted travel prior to the lockdowns in spring 2020.

76

### 77 *Mobility data predicts SARS-CoV-2 spread*

78 We analysed SARS-CoV-2 B.1 (20A) genomes from 10 European countries for which a minimal number of  
79 genomes from the second wave were already available on 3 November, 2020. Using a two-step procedure  
80 that relied on subsampling relative to country-specific case counts (see Methods), we compiled a data set  
81 of close to 4,000 genomes sampled between 29 January and 31 October, 2020 (Extended Data Table 1).  
82 In order to achieve maximum resolution in our evolutionary reconstructions, we constructed a Bayesian  
83 time-measured phylogeographic model that integrates mobility and epidemiological data. Our approach  
84 simultaneously infers phylogenetic history and ancestral movement throughout this history while also  
85 identifying the drivers of spatial spread<sup>10</sup>. We used the latter functionality to determine the most  
86 appropriate mobility or connectivity measure. Specifically, we considered international air transportation  
87 data, the Google COVID-19 Aggregated Mobility Research Dataset (also referred to here as 'mobility data'  
88 for short), as well as Facebook's Social Connectedness Index (SCI), as covariates of phylogeographic spread  
89 (Extended Data Fig. 1). The Google mobility data contains anonymized mobility flows aggregated over  
90 users who have turned on the Location History setting, which is off by default (cfr. Methods). The Social  
91 Connectedness Index reflects the structure of social networks and has been suggested to correlate with

92 the geographic spread of COVID-19 <sup>11</sup>. To help inform the phylogenetic coalescent time distribution, we  
93 parameterized the viral population size trajectories through time as a function of epidemiological case  
94 count data for the countries under investigation.

95

96 Analyses using both time-homogeneous and time-inhomogeneous models offered strong support for  
97 mobility data as a predictor of spatial diffusion whereas air transportation data and SCI offered no  
98 predictive value (Extended Data Table 2). The fact that mobility data encompassing both air and land-  
99 based transport are required to explain COVID-19 spread highlights the need to consider both types of  
100 transport in containment strategies. To ensure that containment strategies were accommodated by our  
101 reconstructions, we further extended our time-inhomogeneous approach to model bi-weekly variation in  
102 the overall rate of spread between countries as a function of mobility (see Methods, Extended Data Table  
103 2).

104

#### 105 *Dynamic viral transmission through time*

106 We use our probabilistic model of spatial spread informed by genomic data, mobility and epidemiological  
107 data to characterize the dynamics of spread throughout the epidemic in Europe. We first focus on the  
108 ratio of introductions over the total viral flow in and out of each country over time and the genetic  
109 structure of country-specific transmission chains (Fig. 1). For the latter, we use a normalized entropy  
110 measure that quantifies the degree of phylogenetic interspersion of country-specific transmission chains  
111 in the SARS-CoV-2 phylogeny (see Methods). Although estimates for individual dispersal between pairs of  
112 countries can also be obtained (Extended Data Fig. 2), we remain cautious in interpreting these as direct  
113 pathways of spread because the genome sampling only covers a restricted set of European countries. The  
114 mobility to and from each country within our 10-country sample covers between 64% and 96% of the  
115 mobility of these countries to/from all countries within Europe (Extended Data Table 3, Extended Data  
116 Figure 3), except for Norway (27%), for which other Scandinavian countries account for considerable  
117 mobility connections (61%), and the UK (49%), for which Ireland accounts for a large fraction of mobility  
118 connections (38%).

119

120 According to the proportion of introductions, we estimate more viral import than export events for  
121 Switzerland, Norway, the Netherlands and Belgium throughout most of the time period under  
122 investigation. According to the estimated phylogenetic entropy, these countries also experienced many  
123 independent transmission chains since the epidemic started to unfold. This is consistent with country-  
124 specific studies; for the first wave in Belgium for example, about 331 individual introductions were  
125 estimated in the ancestry of a limited sample of 740 genomes <sup>12</sup>. For Portugal, we also estimate higher  
126 proportions of introductions early in the first wave but with a subsequent decline to predominantly export  
127 events. France, Italy and Spain on the other hand are characterized by a relatively high viral export during  
128 the first wave. The proportion of introductions remained relatively low for Italy and Spain following the  
129 first wave, while in France these proportions were high from mid-June until the end of July. The absolute  
130 number of transitions in our sample are however low during this time period. These countries also had  
131 comparatively lower entropy values early in the epidemic, with an increase for France by the start of  
132 summer and a more gradual increase over time for Italy. In Spain however, the genetic complexity of  
133 SARS-CoV-2 transmission chains remained limited. In the UK and Germany, the viral flow in and out of the  
134 country was initially relatively balanced. A recent large-scale genomic analysis in the UK indicates that this  
135 can imply very high absolute numbers of cross-country transmissions, as more than 2,800 independent

136 introduction events were identified from the analysis of 26,181 genomes<sup>13</sup>. Although our sample is limited  
137 compared to this UK-focused analysis, our reconstructions also recover major influx from Spain, France  
138 and Italy during the first wave in the UK (Extended Data Fig. 2). We estimate an increase in the proportion  
139 of introductions for the UK from mid-June, indicating an important viral import relative to export around  
140 this time. The phylogenetic entropy also peaked around this time. In Germany, the proportions increased  
141 slightly later in summer with a concomitant rise in phylogenetic entropy.

142

#### 143 *Introductions thrive in low incidence*

144 To assess the impact of summer travel on the second wave in the different countries, we use our genomic-  
145 mobility reconstruction to estimate both the number of lineages persisting in each country and the  
146 number of newly introduced lineages, and how these proliferated early in the second wave. We focus on  
147 a two-month time period between 15 June 2020 – when many EU and Schengen-area countries opened  
148 their borders to other countries – and 15 August, before which the majority of holiday return travel is  
149 expected for many countries. We identify the number of lineages circulating in each country on 15 August,  
150 and determine whether they result from a lineage that persisted since 15 June or from a unique  
151 introduction after this date (independent of the number of descendants for this lineage on 15 August,  
152 Extended Data Fig. 4). In Fig. 2, we plot (1) the ratio of these unique introductions over the total unique  
153 lineages (unique introductions and persisting lineages<sub>1</sub>), (2) the proportion of descendant lineages on  
154 August 15<sup>th</sup> that resulted from the unique introductions over the total descendants circulating on this  
155 date, and (3) the proportion of descendant tips (sampled genomes) after 15 August that resulted from the  
156 unique introductions over the total number of descendant tips (see Methods and Extended Data Fig. 4).  
157 We estimate a posterior mean proportion of unique introductions that is close to or higher than 0.5 except  
158 for Spain and Portugal. This indicates that by 15 August a relatively large fraction of circulating lineages in  
159 each country was spawned by new introductions over summer. Because the B.1.177/20E (EU1) variant  
160 that was predominantly disseminated through summer travel does not appear to be particularly more  
161 transmissible<sup>2</sup>, this was unlikely due to strong intrinsic advantages of the newly introduced viruses.

162

163 The two proportions of descendants from these introductions on 15 August and after this date measure  
164 the relative success of newly introduced lineages compared to persisting lineages, indicating considerable  
165 variation in onward transmission. In Fig. 2, the country estimates are ordered according to decreasing  
166 average incidence during the 15 June – 15 August time period, suggesting that incidence may shape the  
167 outcome of the introductions. In countries that experienced relatively high summer incidence (e.g. Spain,  
168 Portugal, Belgium and France), the introductions lead to comparatively fewer descendants on August 15<sup>th</sup>  
169 or after. We find a significant overall association between incidence and the difference in the logit-scaled  
170 proportion of unique introductions and the logit-scaled proportion of their descendants on August 15<sup>th</sup> ( $P$   
171 = 0.007) as well as between incidence and the difference in the logit-scaled proportion of unique  
172 introductions and the logit-scaled proportion of descendant tips after August 15<sup>th</sup> ( $P$  = 0.019, Extended  
173 Data Figure 5). With comparatively few descendants from introductions (Fig. 2), Norway may to some  
174 extent be an outlier because lineages estimated as persisting in this country could in fact be introductions  
175 from other Scandinavian countries that are not represented in our genome sample. We recover  
176 qualitatively similar, but more variable and statistically unsupported associations between the success of  
177 introductions and incidence for the two-month time periods before and after the 15 June – 15 August  
178 time period (Extended Data Fig. 5). This indicates that the comparatively higher proportion of

179 introductions as well as the more stable and lower incidence between 15 June and 15 August provided  
180 the ideal conditions for a process of genetic drift by which introductions were able to fuel transmission.  
181  
182 Our estimates show that introductions in the UK particularly benefited from the conditions for successful  
183 onward transmission (Fig. 2), with a considerable fraction of introductions originating from Spain  
184 (Extended Data Fig. 6) reflecting the spread of B.1.177/20E (EU1) that rapidly became the most dominant  
185 strain in the UK <sup>2</sup>. Our analysis captures the expansion of this variant as well as that of B.1.160/20A.EU2,  
186 which together account for more than 25% of the genomes in our data set. While Spain was indeed  
187 inferred to be the origin of B.1.177/20E (EU1), the UK also considerably contributed to its spread (Fig. 3).  
188 The earliest introduction from Spain to the UK was estimated around the time Spain opened most EU  
189 borders (21 June, Fig. 3). While introductions from Spain to other countries soon followed, we estimate a  
190 similar rate and amount of spread from the UK to other countries before these other countries also further  
191 disseminated the virus. Although inferred from a limited sample, this illustrates a dynamic pattern of  
192 spread and the importance of the early establishment of B.1.177/20E (EU1) in the UK that likely served as  
193 an important secondary center of dissemination. We note however that this pattern may be impacted by  
194 the intensive and continuous genomic surveillance in the UK, which may also be reflected in our  
195 subsample of the available data. While the UK is also involved in the spread of B.1.160/20A.EU2, this  
196 variant has been largely disseminated from France. The simple fact that this variant expanded later in  
197 France and subsequently also started to spread later compared to B.1.177/20E (EU1) (Extended Data Fig.  
198 7) may explain why the latter spread more successfully.

199 **Discussion**

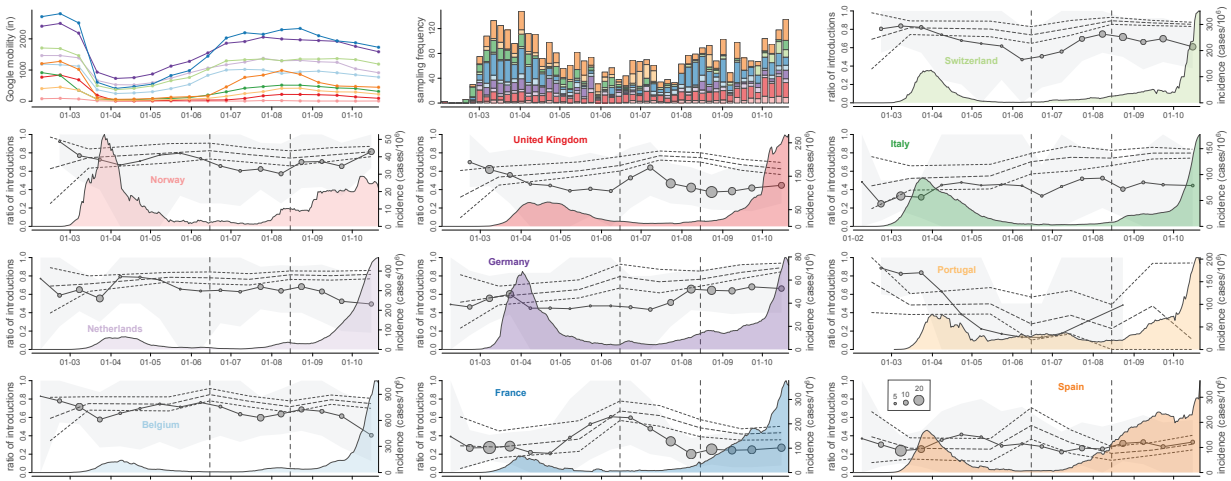
200 Our Bayesian phylogeographic approach builds on a rich history of identifying drivers of spatial spread,  
201 with applications to various pathogens at different spatial scales, ranging from air transportation for  
202 influenza at a global scale <sup>10</sup> to gravity model transmission for Ebola in West Africa <sup>14</sup>. Such studies use a  
203 relatively limited genomic sample to gain insights into viral transmission dynamics. This is also the case in  
204 our application to SARS-CoV-2 in Europe for which we further extend the phylodynamic data integration  
205 approach to confront the lack of resolution offered by SARS-CoV-2 genomic data. A concerted effort in  
206 containing international spread further sets apart the COVID-19 pandemic from these earlier events. For  
207 this reason, we have now incorporated variation in mobility over time to account for the impact of these  
208 measures. Our reconstructions show that the composition of lineages circulating towards the end of the  
209 summer was to an important extent shaped by introductions in most of the European countries. The  
210 relative success of onward transmission of the introduced lineages appears to be shaped by local COVID-  
211 19 incidence during summer.

212  
213 Our results should be interpreted in light of several important limitations. In addition to a limited overall  
214 size, the genome data only cover a selection of European countries, implying that we are missing  
215 transmission events that involve unsampled countries. This may be important for Norway for example,  
216 which according to our mobility data, is largely connected to other Scandinavian countries. We also lack  
217 sampling from eastern Europe, which was to a large extent spared by border controls and lockdowns  
218 during the first wave, but witnessed high excess mortality rates during the second wave. The emergence  
219 of more transmissible variants has led to more intensified genomic surveillance, so similar phylodynamic  
220 reconstructions may now be performed on a wider scale.

221  
222 The pandemic exit strategy offered by vaccination programs is a source of optimism that also sparked  
223 proposals by EU member states to issue vaccine passports in a bid to revive travel and rekindle the  
224 economy. In addition to implementation challenges and issues of fairness, there are risks associated with  
225 such strategies when immunization is incomplete, as likely will be the case for the European population  
226 this summer. A recent modelling study for the United Kingdom suggests that vaccination in adults alone  
227 is unlikely to completely halt the spread of COVID-19 cases and that lifting containment measures early  
228 and suddenly can lead to a large wave of infections <sup>15</sup>. A gradual release of restrictions was shown to be  
229 critical for minimizing the infection burden <sup>15</sup>. We believe that travel policies may be a key consideration  
230 in this respect because similar conditions may arise as the ones we demonstrated to provide fertile ground  
231 for viral dissemination and resurgence in 2020. This may now also involve the spread of variants that are  
232 more transmissible and/or evade immune responses triggered by vaccines and previous infections. Well-  
233 coordinated European strategies will therefore be required to manage the spread of SARS-CoV-2 and  
234 reduce future waves of infection, with hopefully a more unified implementation than hitherto observed.

- 236 1. European Centre for Disease Prevention and Control. Data on 14-day notification rate of new  
237 COVID-19 cases and deaths. [https://www.ecdc.europa.eu/en/publications-data/data-national-14-](https://www.ecdc.europa.eu/en/publications-data/data-national-14-day-notification-rate-covid-19)  
238 [day-notification-rate-covid-19](https://www.ecdc.europa.eu/en/publications-data/data-national-14-day-notification-rate-covid-19) (2021).
- 239 2. Hodcroft, E. B. et al. Spread of a SARS-CoV-2 variant through Europe in the summer of 2020. *Nature*  
240 (2021). <https://doi.org/10.1038/s41586-021-03677-y>.
- 241 3. COVID-19 situation update for the EU/EEA, as of week 3, updated 28 January 2021.  
242 <https://www.ecdc.europa.eu/en/cases-2019-ncov-eueea>.
- 243 4. Rambaut, A., Loman, N., Pybus, O.G., Barclay, W., Barrett, J., Carabelli, A., Connor, T., Peacock, T.,  
244 Robertson, D.L., Volz, E., on behalf of COVID-19 Genomics Consortium UK (CoG-UK). Preliminary  
245 genomic characterisation of an emergent SARS-CoV-2 lineage in the UK defined by a novel set of  
246 spike mutations. [https://virological.org/t/preliminary-genomic-characterisation-of-an-emergent-](https://virological.org/t/preliminary-genomic-characterisation-of-an-emergent-sars-cov-2-lineage-in-the-uk-defined-by-a-novel-set-of-spike-mutations/563)  
247 [sars-cov-2-lineage-in-the-uk-defined-by-a-novel-set-of-spike-mutations/563](https://virological.org/t/preliminary-genomic-characterisation-of-an-emergent-sars-cov-2-lineage-in-the-uk-defined-by-a-novel-set-of-spike-mutations/563) (2020).
- 248 5. Di Domenico, L., Pullano, G., Sabbatini, C. E., Boëlle, P.-Y. & Colizza, V. Impact of lockdown on  
249 COVID-19 epidemic in Île-de-France and possible exit strategies. *BMC Med.* **18**, 1–13 (2020).
- 250 6. Neher, R. A., Dyrda, R., Druelle, V., Hodcroft, E. B. & Albert, J. Potential impact of seasonal forcing  
251 on a SARS-CoV-2 pandemic. *Swiss Med. Wkly* **150**, w20224 (2020).
- 252 7. McKee, M. A European roadmap out of the covid-19 pandemic. *BMJ* **369**, m1556 (2020).
- 253 8. Rambaut, A. et al. A dynamic nomenclature proposal for SARS-CoV-2 lineages to assist genomic  
254 epidemiology. *Nat Microbiol* **5**, 1403–1407 (2020).
- 255 9. Morel, B. et al. Phylogenetic analysis of SARS-CoV-2 data is difficult. *Molecular Biology and*  
256 *Evolution* (2020) doi:10.1093/molbev/msaa314.
- 257 10. Lemey, P. et al. Unifying Viral Genetics and Human Transportation Data to Predict the Global  
258 Transmission Dynamics of Human Influenza H3N2. *PLoS Pathog.* **10**, e1003932 (2014).
- 259 11. Kuchler, T., Russel, D. & Stroebel, J. The geographic spread of COVID-19 correlates with the  
260 structure of social networks as measured by Facebook. NBER Working Paper 26990 (National  
261 Bureau of Economic Research, 2020).
- 262 12. Dellicour, S. et al. A Phylodynamic Workflow to Rapidly Gain Insights into the Dispersal History and  
263 Dynamics of SARS-CoV-2 Lineages. *Mol. Biol. Evol.* (2020) doi:10.1093/molbev/msaa284.
- 264 13. du Plessis, L. et al. Establishment and lineage dynamics of the SARS-CoV-2 epidemic in the UK.  
265 *Science* (2021) doi:10.1126/science.abf2946.
- 266 14. Dudas, G. et al. Virus genomes reveal factors that spread and sustained the Ebola epidemic. *Nature*  
267 **544**, 309–315 (2017).
- 268 15. Moore, S., Hill, E. M., Tildesley, M. J., Dyson, L. & Keeling, M. J. Vaccination and non-pharmaceutical  
269 interventions for COVID-19: a mathematical modelling study. *Lancet Infect. Dis.* (2021)  
270 doi:10.1016/S1473-3099(21)00143-2.

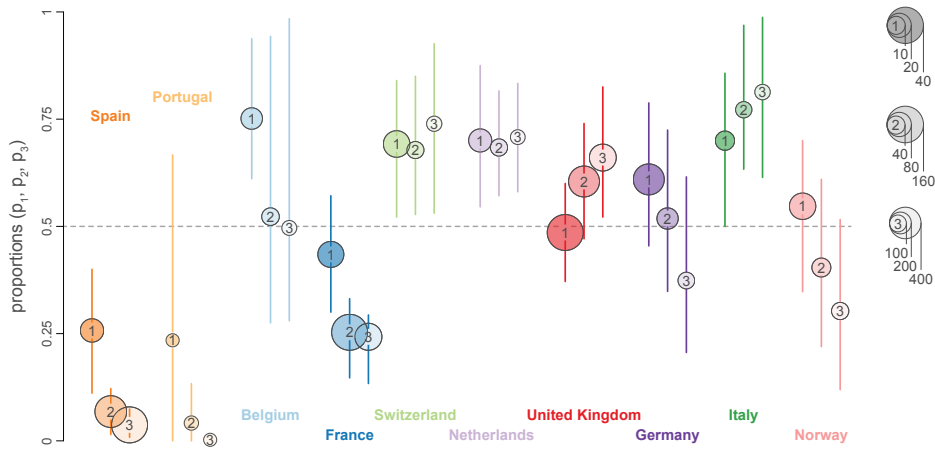




273  
 274

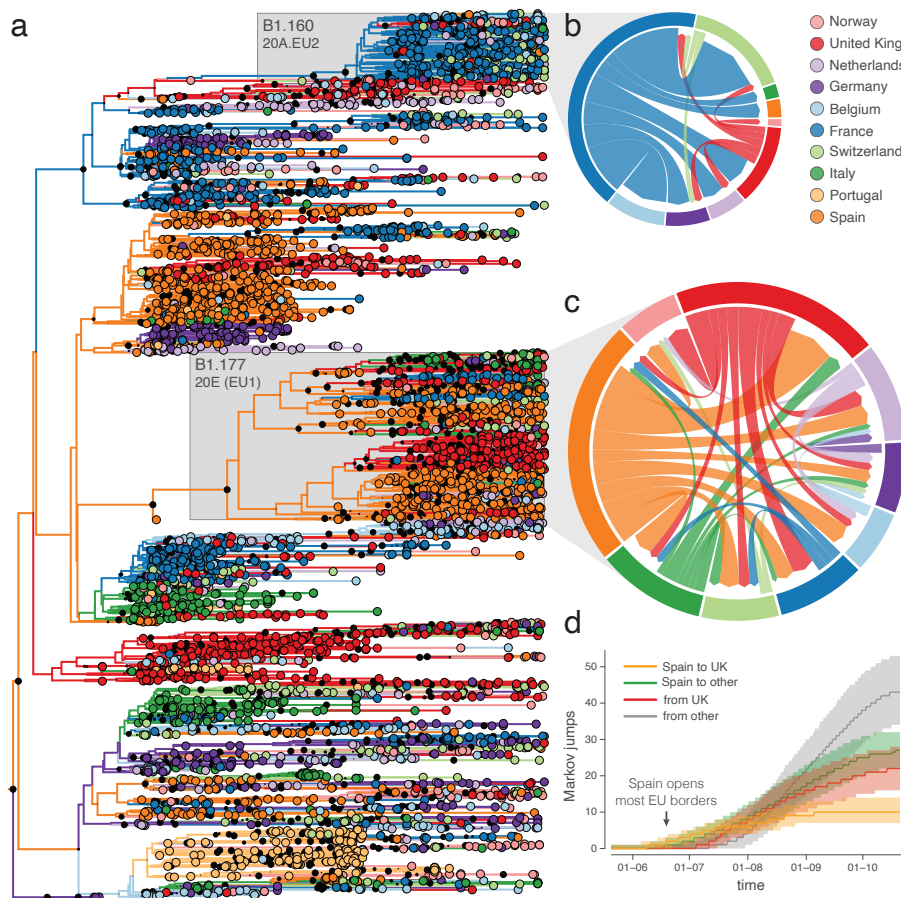
275 **Figure 1. Mobility, genome sampling, case counts and phylogeographic summaries through time for 10 European**  
 276 **countries.** The first panel summarizes the country-specific Google mobility influx in the 10 countries during two-  
 277 week intervals, while the second panel depicts the weekly genome sampling by country used in the phylogeographic  
 278 analysis. In the remaining panels, we plot for each country the ratio of introductions over the total viral flow from  
 279 and to that country (for two-week intervals) and a monthly normalized entropy measure summarizing the  
 280 phylogenetic structure of country-specific transmission chains. The posterior mean ratios of introductions are  
 281 depicted with circles that have a size proportional to the total number of transitions from and to that country and  
 282 the grey surface represents the 95% highest posterior density (HPD) intervals. The posterior mean normalized  
 283 entropies and 95% HPD intervals are depicted by dotted lines. These normalized entropy measures indicate how  
 284 phylogenetically structured the epidemic is in each country, and ranges from 0 (perfectly structured, e.g., a single  
 285 country-specific cluster) to 1 (unstructured interspersed of country-specific sequences across the entire SARS-CoV-  
 286 2 phylogeny). The introduction ratios and normalized entropy measures are superimposed over COVID-19 incidence  
 287 (daily cases/ $10^6$  people) reported for each country through time (coloured density plot). The two vertical dashed  
 288 lines represent the summer time interval (15 June and 15 August, 2020) for which we subsequently evaluate  
 289 introductions versus persistence (see Fig. 2).

290  
 291



292  
 293  
 294  
 295  
 296  
 297  
 298  
 299  
 300  
 301  
 302  
 303  
 304

**Figure 2. Posterior estimates for the relative importance of lineage introduction events in 10 European countries and their association with incidence.** We report three summaries (posterior mean and 95% HPD intervals) for each country: the ratio of unique introductions over the total number of unique persisting lineages and unique introductions between 15 June and 15 August, 2020 ( $p_1$ ), the ratio of descendant lineages from these unique introduction events over the total number of descendants circulating on August 15<sup>th</sup>, 2020 ( $p_2$ ), and the ratio of descendant taxa from these unique introductions over the total number of descendant taxa sampled after 15 August, 2020 ( $p_3$ ) (see Extended Data Fig. 4). The dots are numbered and the sizes are proportional to: (1) the total number of unique lineage introductions identified between 15 June and 15 August, 2020, (2) the total number of lineages inferred on 15 August, 2020, and (3) the total number of descendant tips after 15 August, 2020.



305  
 306 **Figure 3.** Phylogeographic estimates of SARS-CoV-2 spread in 10 European countries. **a**, The maximum clade  
 307 credibility tree summary of the Bayesian inference. Colours correspond to the countries in the legend. The two clades  
 308 corresponding to B1.160/20A.EU2 and B1.177/20E (EU1) are highlighted in grey. **b,c**, Circular migration flow plots  
 309 for for B1.160/20A.EU2 (b) and B1.177/20E (EU1) (c) based on the posterior expectations of the Markov jumps. In  
 310 these plots, migration flow out of a particular location starts close to the outer ring and ends with an arrowhead  
 311 more distant from the destination location. For B1.177/20E (EU1), we also summarize phylogeographic transitions  
 312 as posterior mean estimates with 95% HPD intervals over time for four types of Markov jumps: i) from Spain to the  
 313 UK, ii) from Spain to other countries, iii) from the UK, and iv) from other countries.  
 314

## 315 **Methods**

### 316 *Sequence data and subsampling*

317 We used a two-step genome data collection procedure. We first evaluated the available genomes from  
318 European countries in GISAID <sup>16</sup> on 3 November, 2020. We selected genomes from Belgium, France,  
319 Germany, Italy, Netherlands, Norway, Portugal, Spain, Switzerland and the UK primarily based on the  
320 availability of genome data from both the first and second wave at that time but also because of their  
321 high ratio of genomes to positive cases. A total of 39,812 genomes were available for these countries on  
322 3 November, 2020; the available number of genomes by country are listed in Extended Data Table 1.  
323 Portugal represented an exception because data for this country were limited to the first wave at that  
324 time, but we included genomes from Portugal because of its potential importance as a summer travel  
325 location.

326

327 We aligned the genomes from each country using MAFFT v7.453 <sup>17</sup> and trimmed the 5' and 3' ends and  
328 only retained unique sequences from each location. To further mitigate the disparities in sampling, we  
329 subsampled each country proportionally to the cumulative number of cases on October 21<sup>st</sup> (the most  
330 recently sampled sequence at the time) by setting an arbitrary threshold of 6.5 sequences per 10,000  
331 cases, with a minimum number of 100 sequences per country. To maximize the temporal and spatial  
332 coverage in each country, we binned genomes by epi-week and sampled as evenly as possible, sampling  
333 from a different region within the country when available. Only sequences from the B.1 lineage with the  
334 D614G mutation and exact sampling dates were selected for the analyses. From the final aligned sequence  
335 set, we removed 12 potential outliers, based on a root-to-tip regression applying TempEst v1.5.3 <sup>18</sup> to a  
336 maximum-likelihood tree inferred with IQTREE v2.0.3 <sup>19</sup>, yielding a data set of 2,909 genomes (Extended  
337 Data Table 1).

338

339 Because of the nature of genome sequence accumulation, fewer recently sampled genomes were  
340 available for most countries on 3 November (relative to the case counts at this time). Because our primary  
341 goal was to assess the persistence and introduction of lineages leading up to the second wave, we sought  
342 to augment our data set with more recent genomes, having already performed analyses on the initial data  
343 set. In the section on Bayesian evolutionary reconstructions, we outline how we update these analyses  
344 accordingly. On 5 January, 2021, we updated our dataset by adding over 1,000 non-identical sequences  
345 collected between 1 August and 31 October (out of a total of 56,395 available genomes; the available and  
346 selected number of genomes by country are listed in Extended Data Table 1). For Portugal, we extended  
347 this period back to 22 June (the most recent sampling date for the previous Portuguese selection). We  
348 downloaded all new B.1 sequences with the D614G mutation collected during the selected time period  
349 from GISAID and performed the following subsampling. The number of genomes to add by country was  
350 obtained by raising the threshold ratio of sequences/cases to 8.5 and increasing the minimum number of  
351 sequences to 200. To bias the temporal coverage towards more recent samples, the genomes from each  
352 country were binned by week and sampled such that the number of sequences added by week was  
353 proportional to an exponential function of the form  $e^{t/4}$ , where  $t=0$  represents August 1<sup>st</sup> and  $t=13$  is  
354 October 31<sup>st</sup>. For Portugal, we did not use this preferential sampling as we needed to include close to all  
355 available genomes to raise the number of genomes to 200. The selected sequences were deduplicated  
356 and outliers were removed as described in the previous paragraph. With the additional selection of 1,050  
357 genomes, we arrived at a data set of 3,959 genomes (Extended Data Table 1).

358

359 *Mobility data*

360 We analysed four different mobility and connectivity measures: air traffic flows, a social connectedness  
361 index provided by Facebook, as well as aggregate Facebook<sup>20</sup> and Google international mobility data. Air  
362 traffic flow data were obtained from the International Air Transport Association (<http://www.iata.org>)  
363 and based on the number of origin-destination tickets while also taking into account connections at  
364 intermediate airports<sup>21</sup>. We used monthly air traffic data between the 10 European countries under  
365 investigation for the time period between January 2020 and October 2020. The social connectedness  
366 index (SCI) is an anonymized snapshot of active Facebook users and their friendship networks to measure  
367 the intensity of social connectedness between countries (<https://data.humdata.org/>)<sup>22</sup>. In practice, the  
368 SCI measures the relative probability of a Facebook friendship link between two users of the application  
369 in different countries. We used the SCI calculated for the 10 European countries represented in our  
370 genomic sample as of August 2020.

371

372 The Google COVID-19 Aggregated Mobility Research Dataset contains anonymized mobility flows  
373 aggregated over users who have turned on the Location History setting (on a range of platforms<sup>23</sup>), which  
374 is off by default. To produce this dataset, machine learning is applied to logs data to automatically segment  
375 it into semantic trips<sup>24</sup>. To provide strong privacy guarantees, all trips were anonymized and aggregated  
376 using a differentially private mechanism<sup>25</sup> to aggregate flows over time (see  
377 <https://policies.google.com/technologies/anonymization>). This research was done on the resulting  
378 heavily aggregated and differentially private data. No individual user data was ever manually inspected,  
379 only heavily aggregated flows of large populations were handled. All anonymized trips were processed in  
380 aggregate to extract their origin and destination location and time. For example, if users traveled from  
381 location  $a$  to location  $b$  within time interval  $t$ , the corresponding cell  $(a, b, t)$  in the tensor would be  $n \pm \eta$ ,  
382 where  $\eta$  is Laplacian noise. The automated Laplace mechanism adds random noise drawn from a zero-  
383 mean Laplace distribution and yields  $(\epsilon, \delta)$ -differential privacy guarantee of  $\epsilon = 0.66$  and  $\delta = 2.1 \times 10^{-29}$   
384 per metric. Specifically, for each week  $W$  and each location pair  $(A, B)$ , we compute the number of unique  
385 users who took a trip from location  $A$  to location  $B$  during week  $W$ . To each of these metrics, we add  
386 Laplace noise from a zero-mean distribution of scale  $1/0.66$ . The parameter  $\epsilon$  controls the noise intensity  
387 in terms of its variance, while  $\delta$  represents the deviation from pure  $\epsilon$ -privacy. The closer they are to zero,  
388 the stronger the privacy guarantees. We used aggregated mobility flows between the 10 European  
389 countries and summarized them by two-week or monthly time periods between January 2020 and  
390 October 2020.

391

392 Finally, we also considered international mobility data from Facebook mobility data as an alternative to  
393 Google mobility data. These data are based on numbers of Facebook users moving over large distances,  
394 like air or train travel. Counts of international travel patterns are updated daily based only on users who  
395 have opted to share precise location data from their device with the Facebook mobile app through  
396 location services. Also in this case, we used aggregated mobility flows between the 10 European countries  
397 and summarized them by month between January 2020 and October 2020. Because international  
398 aggregate mobility data obtained from Google and Facebook are highly correlated (monthly Spearman  
399 correlation ranging from 0.84 to 0.92; Supplementary Figure 1), we only included the Google aggregate  
400 mobility data as a covariate in the phylogeographic analyses. We note that the mobility data are subject

401 to limitations as these may not be representative for the population as whole and their representativeness  
402 may vary by location.

403

#### 404 *Bayesian evolutionary reconstructions*

405 - Joint sequence-trait inference with a time-homogeneous generalized linear model of discrete  
406 trait diffusion

407 We performed Bayesian evolutionary reconstruction of timed phylogeographic history using BEAST 1.10  
408 <sup>26</sup> incorporating genome sequences, their country and date of sampling, epidemiological and mobility  
409 and/or connectivity data. Because of the relatively low degree of resolution offered by the sequence data,  
410 our full probabilistic model specification focuses on i) relatively simple model specifications and ii)  
411 informing parameters by additional non-genetic data sources. We modeled sequence evolution using an  
412 HKY85 nucleotide substitution model with gamma-distributed rate variation among sites and a strict  
413 molecular clock model. Our genome set includes three genomes from an early outbreak in Bavaria, which  
414 was caused by an independent introduction from China <sup>27,28</sup>. We therefore constrained these genomes as  
415 an outgroup in the analysis, which according to root-to-tip regression plots as a function of sampling time  
416 resulted in a better correlation coefficient and  $R^2$  compared to the best-fitting root under the heuristic  
417 mean residual squared criterion (Supplementary Figure 2) <sup>18</sup>.

418

419 As a coalescent tree prior, we modeled the effective population size trajectory as a piecewise constant  
420 function that changes values at pre-specified times (following <sup>29</sup>), with log population sizes modelled as a  
421 deterministic function of log COVID-19 case counts (following <sup>30</sup>). This reduces the nonparametric skygrid  
422 parameterization to a generalized linear model (GLM) formulation with an estimable regression intercept  
423 ( $\alpha$ ) and coefficient ( $\beta$ ). In this parameterization, a coefficient estimate centered around 0 would imply  
424 constant population size dynamics through time. We specified two-week intervals and summarized as a  
425 covariate the total case counts over these time intervals for the 10 countries of sampling (obtained from  
426 <https://www.ecdc.europa.eu/en/covid-19/data>). The earliest interval with non-zero cases counts was  
427 from 2020-01-14 to 2020-01-28; before 2020-01-14, the log-transformed and standardized case count  
428 covariate was set to the equivalent of 1 case. We also tested whether a lag-time was required for the case  
429 count covariate using marginal likelihood estimation (MLE) <sup>31</sup>. Specifically, we shifted the case counts by  
430 1, 2, 3 and 4 weeks before summarizing them according to two-week intervals and estimated the model  
431 fit of these covariates against case counts without lag time (Supplementary Table 1). To mitigate the  
432 computational burden associated with the MLE procedure, we performed these analyses on a subset of  
433 1,000 genomes (obtained using the Phylogenetic Diversity Analyzer tool <sup>32</sup>). We estimated the highest  
434 (log) marginal likelihood for a two-week lag time (Supplementary Table 1) and used this for the case count  
435 covariate in our analyses.

436

437 Similar to sequence evolution, we modelled the process of transitioning through discrete location states  
438 (countries of sampling) according to a continuous-time Markov chain (CTMC) <sup>33</sup>. We employed a  
439 parameterization that models the log transition rates as a log linear function of mobility and connectivity  
440 covariates <sup>10</sup>. The Bayesian implementation of this model simultaneously estimates phylogenetic history,  
441 ancestral movement and the contribution of covariates to the movement patterns <sup>10</sup>. While we mainly  
442 use this approach to obtain well-informed phylodynamic estimates, we also make use of its capacity to  
443 identify the most relevant mobility measure to inform our reconstructions. As covariates we considered

444 Facebook's SCI, air transportation data and mobility data. For the two time-variable mobility measures,  
445 we used the average of the log-transformed and standardized monthly mobility measures as a single  
446 covariate in our time-homogeneous phylogeographic GLM model. In this GLM formulation, we estimate  
447 positive effect sizes for each covariate as well as their inclusion probability through a spike-and-slab  
448 procedure<sup>10</sup>. Although we subsampled the number of SARS-CoV-2 genomes by country in proportion to  
449 case counts, they do not fully correspond because we used a minimum number of genomes for countries  
450 with low case counts. We therefore evaluated whether this resulted in signal for sampling bias by including  
451 an origin and destination covariate in the GLM based on the residuals for a regression analysis between  
452 genomes and case counts (following<sup>14</sup>). We performed this analysis using a set of empirical trees (see  
453 Time-inhomogeneous reconstructions) applying both a time-homogeneous and time-inhomogeneous  
454 model, but found no support for these additional covariates (Supplementary Table 2).

455  
456 We performed inference under the full model specification using Markov chain Monte Carlo (MCMC)  
457 sampling and used the BEAGLE library v3<sup>34</sup> to increase computational performance. We specified  
458 standard transition kernels on all parameters, except for the regression coefficients of the piecewise-  
459 constant coalescent GLM model. For these parameters, we implemented new Hamiltonian Monte Carlo  
460 (HMC) transition kernels to improve sampling efficiency. These kernels use principles from Hamiltonian  
461 dynamics and their approximate energy conserving properties to reduce correlation between successive  
462 sampled states, but require computation of the gradient of the model log-posterior with respect to the  
463 parameters of interest, in addition to efficient evaluation of the log-posterior that BEAGLE provides. To  
464 accomplish this, we extended our previous analytic derivation of the gradient of the log-density from the  
465 skygrid coalescent model with respect to the log-population-sizes<sup>35</sup> to now be with respect to the  
466 regression coefficients using the chain rule and their regression design matrix.

467  
468 Due to the data set size, MCMC burn-in takes up considerable computational time. We therefore iterated  
469 through a series of BEAST inferences, initially only considering sequence evolution and subsequently  
470 adding the location data, to arrive at a tree distribution from which trees were taken as starting trees in  
471 our final analyses. The latter was composed of multiple independent MCMC runs that were run sufficiently  
472 long to ensure that their combined posterior samples achieved effective sample sizes (ESSs) larger than  
473 100 for all continuous parameters.

474  
475 - Data augmentation through online BEAST

476 As we updated our dataset following initial analyses of the 2,909 genome collection using the approach  
477 discussed (see Bayesian evolutionary reconstructions), we sought to capitalize on these efforts to limit  
478 the burn-in for subsequent analyses of the 3,959 dataset. Specifically, we adopted the distance-based  
479 procedure to insert new taxa into a time-measured phylogenetic tree sample as implemented in the  
480 BEAST framework for online inference<sup>36</sup>. We subsequently use the augmented tree as the starting tree  
481 for the analyses of the updated dataset.

482  
483 - Time-inhomogeneous reconstructions

484 To accommodate the time-variability of the mobility measures, we constructed epoch model extensions  
485 of the discrete phylogeography approach that allow specifying arbitrary intervals over the evolutionary  
486 history and associating them with different model parameterizations<sup>37</sup>. As a complement to testing

487 covariates of spatial diffusion using a time-homogeneous model, we used the epoch extension to specify  
488 monthly intervals allowing us to incorporate monthly mobility matrices (air transportation data were only  
489 available as monthly numbers), but assuming time-homogeneous effect sizes and inclusion probabilities.  
490 Monthly covariates were again log-transformed and standardized after adding a pseudo-count to each  
491 entry in the monthly matrices.

492

493 In addition, we performed another analysis in which we relaxed the constant-through-time inclusion  
494 probability of the covariates. In this model specification, each interval is associated with a specific set of  
495 indicator variables to represent the inclusion/exclusion of covariates, but we pool information about  
496 predictor inclusion across the intervals using hierarchical graph modelling<sup>38</sup>. This approach uses a set of  
497 indicator variables to model covariate inclusion at the hierarchical level but allows interval-specific  
498 inclusion or predictors to diverge from the hierarchical level with a non-zero probability (with the number  
499 of differences modelled as a binomial distribution<sup>38</sup>), which was set to 0.10 in our case. We estimated  
500 hierarchical and interval-level inclusion using spike-and-slab<sup>38</sup>.

501

502 Finally, we performed an analysis using the time-inhomogeneous model in which the interval-specific  
503 transition rates are modelled as a function of the single covariate that is supported by the analyses above  
504 leveraging aggregate mobility. We incorporated more variability through time by specifying two-week  
505 intervals (similar to the coalescent GLM interval specification). In addition, we add time-homogeneous  
506 random effects to the phylogeographic transition rate parameterization in order to account for potential  
507 biases in the ability of mobility to predict phylogeographic spread. While posterior mean estimates for  
508 these random effects vary, only very few indicate that individual phylogeographic transition rates  
509 significantly deviate from the mobility data (Supplementary Figure 3). The time-inhomogeneous GLM  
510 approach we employ allows modelling relative differences in transition rates, but also the overall rate of  
511 migration between countries varies through time, and importantly, this is strongly affected by  
512 intervention strategies. To accommodate these dynamics, we further extended this model by  
513 incorporating a time-inhomogeneous overall CTMC rate scaler and parameterize it as a log linear function  
514 of the total monthly between-country log-transformed and standardized mobility (time-variable rate  
515 scalar GLM in Extended Data Table 2). To generate realisations of the discrete location CTMC process and  
516 obtain estimates of the transitions (Markov jumps) between states under this model, we employed  
517 posterior inference of the complete Markov jump history through time<sup>10,39</sup>.

518

519 While the epoch model allows us to flexibly accommodate time-variable spatial dynamics, it considerably  
520 increases the computational burden associated with likelihood evaluations. In order to efficiently draw  
521 inference under this model for our large data set, we fit the time-inhomogeneous spatial diffusion process  
522 to a set of trees inferred under the time-homogeneous GLM diffusion model described above. Although  
523 likelihood evaluations remain computationally expensive, even with the speed-up offered by GPU  
524 computation with BEAGLE, eliminating simultaneous tree estimation tremendously reduces parameter-  
525 space, requiring only modest MCMC chain lengths to adequately explore it. Model and inference  
526 specifications for the different analyses are available as BEAST XML input files on GitHub  
527 ([https://github.com/phylogeography/SARS-CoV-2\\_EUR\\_PHYLOGEOGRAPHY](https://github.com/phylogeography/SARS-CoV-2_EUR_PHYLOGEOGRAPHY)) and Zenodo  
528 (<https://doi.org/10.5281/zenodo.4876442>).

529

530 - Posterior Summaries



531 We assessed MCMC mixing (e.g. using ESSs) and summarized continuous parameter estimates using  
532 Tracer v1.7.1<sup>40</sup>. Credible intervals were computed as 95% HPD intervals. Trees were visualized using  
533 FigTree v1.4.4 (available at <https://github.com/rambaut/figtree/releases>). In terms of phylogeographic  
534 estimates, we mainly focused on i) transitions to each location and from each location (based on Markov  
535 jump estimates) instead of pairwise transitions, ii) ratios of these transitions and iii) how these transitions  
536 structured transmission chains in individual countries. Transitions to each and from each location avoid  
537 drawing conclusions about direct migration between countries, which can be tenuous given the  
538 incomplete genomes coverage of Europe, while their ratios avoid using absolute numbers of transitions,  
539 which are highly sample-dependent. Phylogeographic inference is limited to reconstructing the transitions  
540 in the ancestral history of a sample of sequences, which will only be a small fraction of the actual migration  
541 events especially when these events result in insufficient onward transmission to be captured in our  
542 limited sample. In addition, SARS-CoV-2 genome data can be poorly resolved and identical genomes in  
543 different locations are consistent with hypotheses that involve both a sparse and a rich number of virus  
544 flows between these locations. As the data hold little information to distinguish these hypotheses, we  
545 only consider sparse scenario's by including only unique sequences for each location. A joint inference of  
546 sequence evolution and discrete spatial diffusion would err on the side of sparse hypotheses anyway  
547 because it will tend to cluster identical sequences that share a location. Despite the general  
548 underestimation of spatial dispersal, a phylogeographic inference is still likely to capture the transition  
549 events with important onward transmission, and evaluating the importance of such events relative to  
550 persistence is a major focus of this study. Cryptic transmission also complicates the ability to reconstruct  
551 spatial dispersal, but we expect this to be equally likely for introductions and persistence and therefore  
552 focus on their ratio for each location.

553

554 We provide three new tree sample tools in the BEAST codebase available at [https://github.com/beast-](https://github.com/beast-dev/beast-mcmc)  
555 [dev/beast-mcmc](https://github.com/beast-dev/beast-mcmc)) to obtain posterior summaries of location transition histories using posterior tree  
556 distributions annotated with Markov jumps:

557

558 • *TreeMarkovJumpHistoryAnalyzer* allows collecting Markov jumps and their timings from a  
559 posterior tree distribution annotated with Markov jumps histories in a .csv file for further  
560 analyses.

561

562 • *TreeStateTimeSummarizer* decomposes the total tree time into the times associated with  
563 contiguous partitions of a tree estimated to be in a particular location state, with the partitions  
564 determined by the Markov jumps. An arbitrary lower- and upper-time boundary can be specified  
565 to restrict the summary to a particular time interval in the evolutionary history. We use the time  
566 estimates for the separate partitions associated with each state to calculate an entropy measure  
567 that summarizes the genetic make-up of country-specific transmission chains. Specifically, we use  
568 for each location a normalized Shannon entropy:

569

$$-\frac{1}{\ln(n)} \sum_i^n p_i \ln(p_i), \quad (1)$$

570 where  $p_i$  is the proportion of time associated with that location for partition  $i$  of a phylogeographic  
571 tree and  $n$  represents the number of partitions for that location in the tree.

572

573 • *PersistenceSummarizer* also uses posterior tree distributions annotated with Markov jumps to  
574 summarize the number of lineages at a particular point in time (evaluation time,  $T_e$ , see Extended

575 Figure 5), which location states they are associated with, since what time point in the past they  
576 have maintained that state and how many sampled descendants they have after time  $T_e$   
577 (Extended Figure 5). In addition, it allows summarizing how long these lineages have circulated  
578 independently prior to  $T_e$ , so before sharing common ancestry with other lineages that  
579 maintained the same location state. This information allows us to determine how many lineages  
580 are circulating at  $T_e$  that stem either from a unique persistent lineage (maintaining the same  
581 location states) or unique introduction event since a particular time prior to  $T_e$  ( $T_a$  in Extended  
582 Figure 5). The association between incidence and the difference in the logit proportion of unique  
583 introductions and the logit proportion of their descendants on August 15<sup>th</sup> was evaluated using a  
584  $p$ -value obtained by a linear regression analysis.

585

## 586 **Data availability**

587 BEAST XML input files are available at  
588 [https://github.com/phylogeography/SARS-CoV-2\\_EUR\\_PHYLOGEOGRAPHY](https://github.com/phylogeography/SARS-CoV-2_EUR_PHYLOGEOGRAPHY) (DOI:  
589 10.5281/zenodo.4876442). The SARS-CoV-2 genome data required for running these XML files can be  
590 downloaded from <https://www.gisaid.org>; all GISAID accession numbers are listed in the GISAID  
591 acknowledgments table (Supplementary Table 3).  
592 The Google COVID-19 Aggregated Mobility Research Dataset and the Facebook mobility data are not  
593 publicly available owing to stringent licensing agreements. Information on the process of requesting  
594 access to the Google mobility data is available from A.S. ([sadilekadam@google.com](mailto:sadilekadam@google.com)) and the COVID-19  
595 Community Mobility Reports that were derived from the Google data are publicly available at  
596 <https://www.google.com/covid19/mobility/>. The Facebook mobility data are made available through  
597 the Data for Good program (<https://dataforgood.fb.com>) under the terms of a data license agreement  
598 which defines the allowed terms of use by partners (contact: [disastermaps@fb.com](mailto:disastermaps@fb.com)). Once a partner  
599 institution's request for access is vetted and an appropriate data license agreement is signed, then  
600 access is granted through a Facebook's web-based spatial visualization tool called GeoInsight. Air travel  
601 data were obtained from the International Air Transport Association (<http://www.iata.org>).  
602 Log-transformed and standardized among country mobility and air travel data are specified in the  
603 available XML files. COVID-19 incidence data was obtained from [https://www.ecdc.europa.eu/en/covid-](https://www.ecdc.europa.eu/en/covid-19/data)  
604 [19/data](https://www.ecdc.europa.eu/en/covid-19/data).

605

## 606 **Code availability**

607 The code for running BEAST analyses is available in the hmc\_develop branch of the BEAST codebase  
608 available at <https://github.com/beast-dev/beast-mcmc> (DOI: 10.5281/zenodo.4895235). The tools  
609 *TreeMarkovJumpHistoryAnalyzer*, *TreeStateTimeSummarizer* and *PersistenceSummarizer* are available  
610 from the master branch in the same codebase.

611

- 612 16. Shu, Y. & McCauley, J. GISAID: Global initiative on sharing all influenza data - from vision to reality.  
613 *Euro Surveill.* **22**, (2017).
- 614 17. Katoh, K., Asimenos, G. & Toh, H. Multiple alignment of DNA sequences with MAFFT. *Methods Mol.*  
615 *Biol.* **537**, 39–64 (2009).
- 616 18. Rambaut, A., Lam, T. T., Max Carvalho, L. & Pybus, O. G. Exploring the temporal structure of  
617 heterochronous sequences using TempEst (formerly Path-O-Gen). *Virus Evol* **2**, vew007 (2016).

- 618 19. Minh, B. Q. *et al.* IQ-TREE 2: New Models and Efficient Methods for Phylogenetic Inference in the  
619 Genomic Era. *Mol. Biol. Evol.* **37**, 1530–1534 (2020).
- 620 20. Maas, P. Facebook disaster maps: Aggregate insights for crisis response & recovery. in *Proceedings*  
621 *of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining* (ACM,  
622 2019). doi:10.1145/3292500.3340412.
- 623 21. Gilbert, M. *et al.* Preparedness and vulnerability of African countries against importations of COVID-  
624 19: a modelling study. *Lancet* **395**, 871–877 (2020).
- 625 22. Bailey, M., Cao, R., Kuchler, T., Stroebel, J. & Wong, A. Social Connectedness: Measurement,  
626 Determinants, and Effects. *J. Econ. Perspect.* **32**, 259–280 (2018).
- 627 23. Kraemer, M. U. G. *et al.* Mapping global variation in human mobility. *Nat Hum Behav* **4**, 800–810  
628 (2020).
- 629 24. Bassolas, A. *et al.* Hierarchical organization of urban mobility and its connection with city livability.  
630 *Nat. Commun.* **10**, 4817 (2019).
- 631 25. Wilson, R. J. *et al.* Differentially private SQL with bounded user contribution. In *Proc. Privacy*  
632 *Enhancing Technologies* 230–250 (2020).
- 633 26. Suchard, M. A. *et al.* Bayesian phylogenetic and phylodynamic data integration using BEAST 1.10.  
634 *Virus Evol* **4**, vey016 (2018).
- 635 27. Böhmer, M. M. *et al.* Investigation of a COVID-19 outbreak in Germany resulting from a single  
636 travel-associated primary case: a case series. *Lancet Infect. Dis.* **20**, 920–928 (2020).
- 637 28. Worobey, M. *et al.* The emergence of SARS-CoV-2 in Europe and North America. *Science* **370**, 564–  
638 570 (2020).
- 639 29. Gill, M. S. *et al.* Improving Bayesian population dynamics inference: a coalescent-based model for  
640 multiple loci. *Mol. Biol. Evol.* **30**, 713–724 (2013).
- 641 30. Faria, N. R. *et al.* Genomic and epidemiological monitoring of yellow fever virus transmission  
642 potential. *Science* **361**, 894–899 (2018).
- 643 31. Baele, G., Lemey, P. & Suchard, M. A. Genealogical Working Distributions for Bayesian Model  
644 Testing with Phylogenetic Uncertainty. *Systematic Biology* vol. 65 250–264 (2016).
- 645 32. Chernomor, O. *et al.* Split diversity in constrained conservation prioritization using integer linear  
646 programming. *Methods Ecol. Evol.* **6**, 83–91 (2015).
- 647 33. Lemey, P., Rambaut, A., Drummond, A. J. & Suchard, M. A. Bayesian phylogeography finds its roots.  
648 *PLoS Comput. Biol.* **5**, e1000520 (2009).
- 649 34. Ayres, D. L. *et al.* BEAGLE 3: Improved performance, scaling, and usability for a high-performance  
650 computing library for statistical phylogenetics. *Syst. Biol.* **68**, 1052–1061 (2019).
- 651 35. Baele, G., Gill, M. S., Lemey, P. & Suchard, M. A. Hamiltonian Monte Carlo sampling to estimate  
652 past population dynamics using the skygrid coalescent model in a Bayesian phylogenetics  
653 framework. *Wellcome Open Res* **5**, 53 (2020).
- 654 36. Gill, M. S., Lemey, P., Suchard, M. A., Rambaut, A. & Baele, G. Online Bayesian Phylodynamic  
655 Inference in BEAST with Application to Epidemic Reconstruction. *Mol. Biol. Evol.* **37**, 1832–1842  
656 (2020).
- 657 37. Bielejec, F., Lemey, P., Baele, G., Rambaut, A. & Suchard, M. A. Inferring heterogeneous  
658 evolutionary processes through time: from sequence substitution to phylogeography. *Syst. Biol.* **63**,  
659 493–504 (2014).
- 660 38. Cybis, G. B., Sinsheimer, J. S., Lemey, P. & Suchard, M. A. Graph hierarchies for phylogeography.  
661 *Philos. Trans. R. Soc. Lond. B Biol. Sci.* **368**, 20120206 (2013).

- 662 39. Minin, V. N. & Suchard, M. A. Fast, accurate and simulation-free stochastic mapping. *Philos. Trans.*  
663 *R. Soc. Lond. B Biol. Sci.* **363**, 2985–2995 (2008).  
664 40. Rambaut, A., Drummond, A. J., Xie, D., Baele, G. & Suchard, M. A. Posterior Summarization in  
665 Bayesian Phylogenetics Using Tracer 1.7. *Syst. Biol.* **67**, 901–904 (2018).

666

## 667 **Acknowledgments**

668 We would like to thank all the authors who have kindly shared genome data on GISAID, and we have  
669 included a table (Supplementary Table 3) acknowledging the authors and institutes involved.

670

671 The research leading to these results has received funding from the European Research Council under the  
672 European Union's Horizon 2020 research and innovation programme (grant agreement no. 725422-  
673 ReservoirDOCS) and the Bill & Melinda Gates Foundation (OPP1094793 and INV-024911). This study was  
674 partially funded by EU grant 874850 MOOD and is catalogued as MOOD 005. The contents of this  
675 publication are the sole responsibility of the authors and do not necessarily reflect the views of the  
676 European Commission. The Artic Network receives funding from the Wellcome Trust through project  
677 206298/Z/17/Z. PL acknowledges support by the Research Foundation - Flanders ('Fonds voor  
678 Wetenschappelijk Onderzoek - Vlaanderen', G066215N, G0D5117N and G0B9317N). GB acknowledges  
679 support from the 'Interne Fondsen KU Leuven' / Internal Funds KU Leuven under grant agreement  
680 C14/18/094, and the Research Foundation – Flanders ('Fonds voor Wetenschappelijk Onderzoek -  
681 Vlaanderen', G0E1420N, G098321N). MAS acknowledges support from National Institutes of Health grant  
682 U19 AI135995 and R01 AI153044. SD is supported by the *Fonds National de la Recherche Scientifique*  
683 (FNRS, Belgium). We gratefully acknowledge support from NVIDIA Corporation with the donation of  
684 parallel computing resources used for this research. We would also like to thank AMD for the donation of  
685 critical hardware and support resources from its HPC Fund that made this work possible.

686

## 687 **Author contributions**

688 P.L. and S.D. designed the study, performed analyses and drafted the manuscript. V.C., C.P. and A.S.  
689 provided and analysed data. S.L.H., F.V.d.B., N.R., S.L. and A.J.T. compiled and analysed data. A.L., B.B.O.M.  
690 and M.K. contributed data. G.B. performed data analyses. M.S.G., X.J. and M.A.S. developed statistical  
691 inference methodology. All authors contributed to interpreting and reviewing the manuscript.

692

## 693 **Competing Interests**

694 The authors declare no competing interests.

695

## 696 **Materials and correspondence**

697 [philippe.lemey@kuleuven.be](mailto:philippe.lemey@kuleuven.be) & [simon.dellicour@ulb.ac.be](mailto:simon.dellicour@ulb.ac.be)

698

699 **Extended Data Figures**

700

701 **Extended Data Figure 1.** Monthly international mobility data matrices: international air traffic data (a), international  
702 Facebook mobility data (b), and international Google mobility data (c). For Facebook data, we also report the single  
703 social connectedness index matrix (SCI, b).

704

705 **Extended Data Figure 2.** Estimated introductions through time in the 10 European countries and circular migration  
706 flow plots summarizing the estimated transitions between the countries for different time intervals throughout the  
707 SARS-CoV-2 evolutionary history. (a) The introductions through time serve as an illustration and are based on the  
708 Markov jump history in the MCC tree. We note that the posterior distribution of trees is accompanied with  
709 considerable uncertainty about the location of origin, destination and timing of the transitions, which is difficult to  
710 appropriately visualize. The grey box represents the time period from 15 June to 15 August. (b) The circular migration  
711 flow plots are based on the posterior expectations of the Markov jumps. The sizes of the plots reflect the total  
712 number of transitions for each period. In these plots, migration flow out of a particular location starts close to the  
713 outer ring and ends with an arrowhead more distant from the destination location.

714

715 **Extended Data Figure 3.** Pairwise mobility data among the 10 countries included in the phylogeographic analysis  
716 and other European countries. Heatmap cells are coloured according to international Google mobility data for the  
717 time period between January and October 2020.

718

719 **Extended Data Figure 4.** Conceptual representation of persistent lineages and introductions during the time interval  
720 delineated by the evaluation time ( $T_e$ ) and the ancestral time ( $T_a$ ). At  $T_e$ , we evaluate how many lineages are  
721 circulating in the location of interest, in this case 12 (lineages in other locations are represented by thick grey  
722 branches). We subsequently identify whether these lineages maintained this location up to  $T_a$  in their ancestry or  
723 whether they result from an introduction event in the time interval of interest. By determining whether other  
724 lineages circulating in the location of interest at  $T_e$  are descendants of the same persistent lineage or whether they  
725 share an introduction event, we identify the unique persistent lineages or introductions, in this case 2 and 4  
726 respectively. In addition to the proportion of unique introductions ( $p_1 = 4/6$ ), we also summarize the proportion of  
727 their descendants at  $T_e$  ( $p_2 = 9/(9+3)$  in this case) and the proportion of their descendants in terms of sampled tips  
728 after  $T_e$  ( $p_3$ ). Those tips are not shown here but conceptually represented for both introductions and persistent  
729 lineages by ovals.

730

731 **Extended Data Figure 5.** Scatter plots of the difference in the logit proportion of unique introductions ( $p_1$ ) and the  
732 logit proportion of their descendants on 15 August ( $p_2$ ) against incidence and the difference in the logit proportion  
733 of unique introductions and the logit proportion of descendant tips after 15 August ( $p_3$ ) against incidence. Both plots  
734 are shown for the period between 15 April and 15 June, for the period between 15 June and 15 August, and for the  
735 period between 15 August and 15 October, respectively. The p-values in the lower right corner of the plots are the  
736 p-values for the hypothesis tests based on the t-statistic evaluating whether the regression coefficient in a linear  
737 regression model is different from 0.

738

739 **Extended Data Figure 6.** Estimated geographic origin of viral influx over the summer (15 June – 15 August, 2020) in  
740 each country. Each bar plot summarizes the posterior Markov jump estimates into a specific country. For the bar  
741 representing a low number of introductions into Portugal, a magnified view is provided.

742

743 **Extended Data Figure 7.** Phylogeographic transitions for lineages B1.1777/20A.EU1 and B1.160/20A.EU2.  
744 Cumulative phylogeographic transitions are summarized as posterior mean estimates with 95% HPD intervals over  
745 time for four types of Markov jumps. For B1.1777/20A.EU1: i) from Spain to the UK, ii) from Spain to other  
746 countries, iii) from the UK, and iv) from other countries; For B1.160/20A.EU2: i) from France to the UK, ii) from  
747 France to other countries, iii) from the UK, and iv) from other countries.

748 **Extended Data Tables**

749

750 **Extended Data Table 1.** Genome sampling by country, collected on November 3<sup>rd</sup>, 2020, and updated on January  
751 5<sup>th</sup>, 2021.

752

753 The numbers in between brackets represent the number of available genomes that were subsampled. \*For Portugal,  
754 almost all available genomes were included to increase the number of genomes to 200.

755

756 **Extended Data Table 2.** Parameter estimates for the various Bayesian time-measured phylogeographical models.

757

758 The coalescent generalized linear model (GLM) parameterizes bi-weekly effective population sizes as a log-linear  
759 function of COVID-19 incidence data, with  $\alpha$  and  $\beta$  representing the log intercept and log regression coefficient. In  
760 the time-inhomogeneous spatial diffusion models, no coalescent prior was used as these models were fitted onto  
761 posterior trees inferred from the time-homogeneous model (see Methods). For the spatial GLM model, we report  
762 inclusion probability estimates through the expectations of the boolean indicators ( $\delta$ ) associated with each predictor  
763 and log conditional effect sizes (the regression coefficient conditional on the predictor being included in the model,  
764  $\beta|\delta=1$ ). SCI = Social Connectedness Index, based on Facebook data. For the model with time-variable inclusion  
765 probabilities, we report the parameters at the hierarchical level ( $\delta_i$  and  $\beta|\delta_i$ , see Methods). In the model with a time-  
766 variable rate scalar, we parameterize this rate scalar as a log-linear function of the overall between-country mobility,  
767 with  $\alpha$  and  $\beta$  representing the log intercept and log regression coefficient.

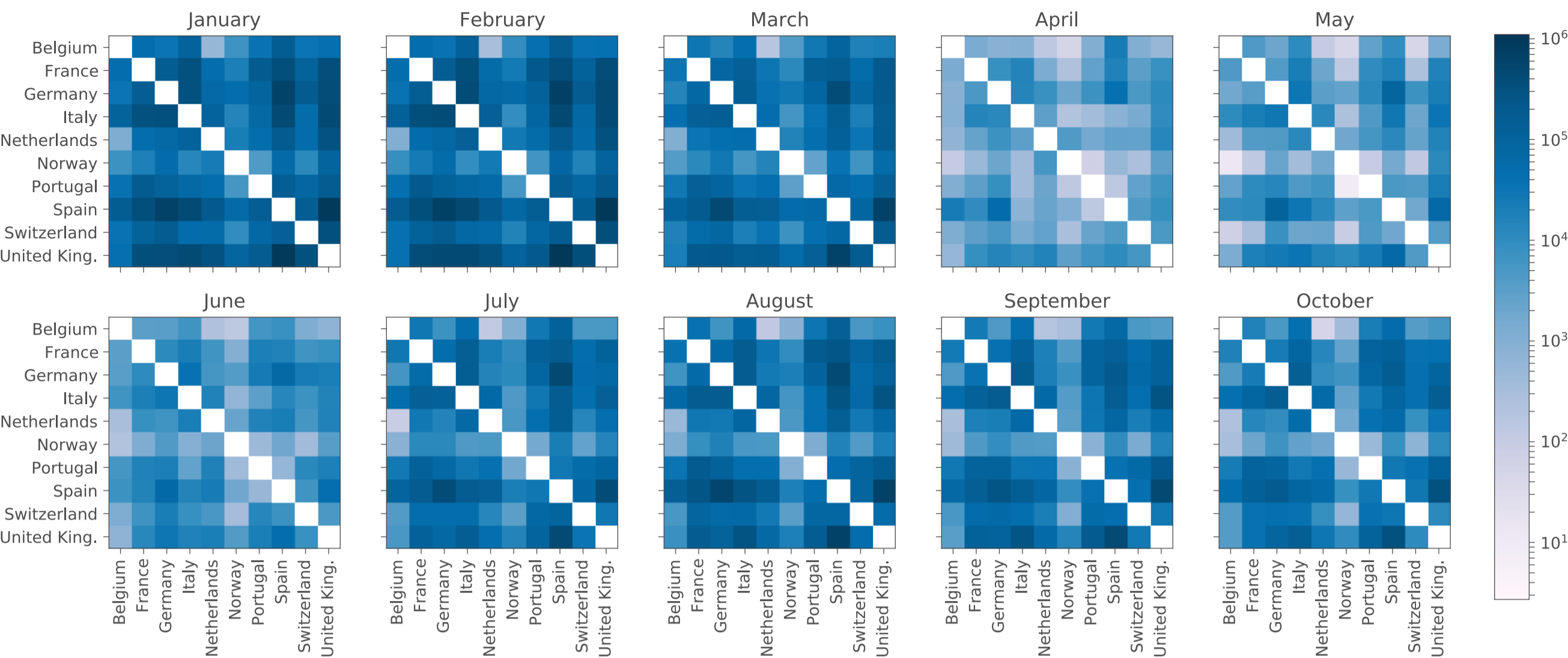
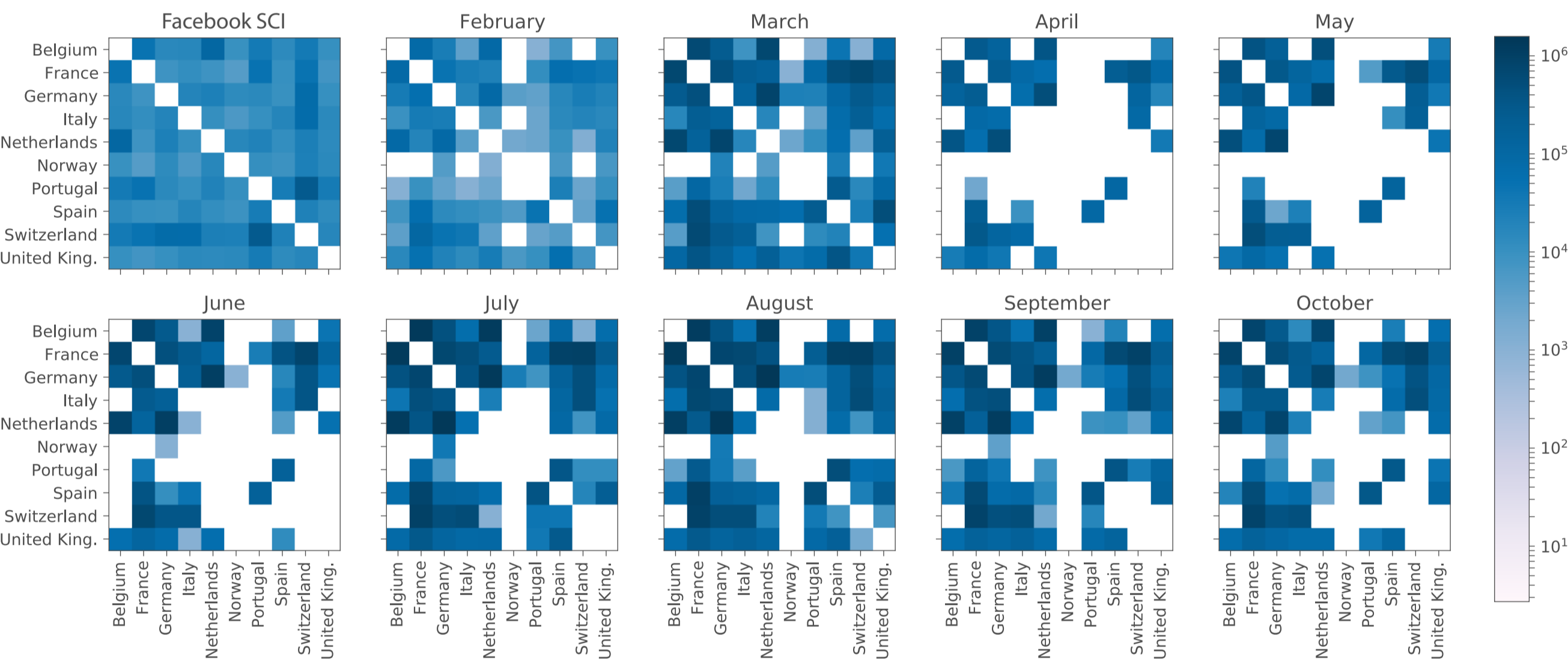
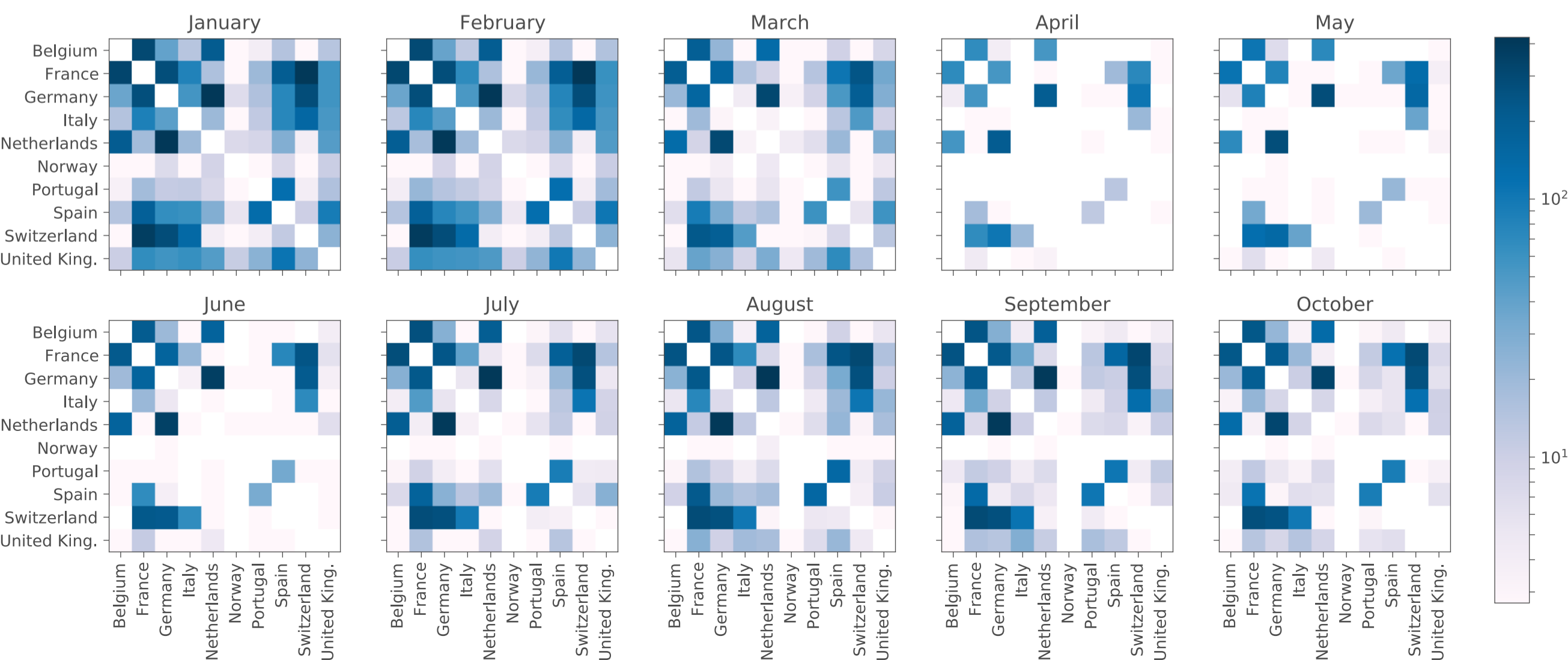
768 Using a time-homogeneous model of spatial diffusion, we estimate a maximum inclusion probability for the mobility  
769 data whereas air transportation data and SCI offer no predictive value. We also estimate a strong positive association  
770 between viral population size change through time and COVID-19 incidence in the coalescent GLM. We further  
771 confirm the support for the mobility covariate in a time-inhomogeneous spatial model that incorporates monthly  
772 mobility measures, with either constant or time-variable inclusion probabilities. In addition to parameterizing the  
773 relative rates of spread between countries according to this covariate, we extend our time-inhomogeneous  
774 approach to also model bi-weekly variation in the overall rate of spread between countries as a function of mobility  
775 measures (time-variable rate scalar GLM). This approach estimates a positive association between the overall rate  
776 of spatial spread and mobility data.

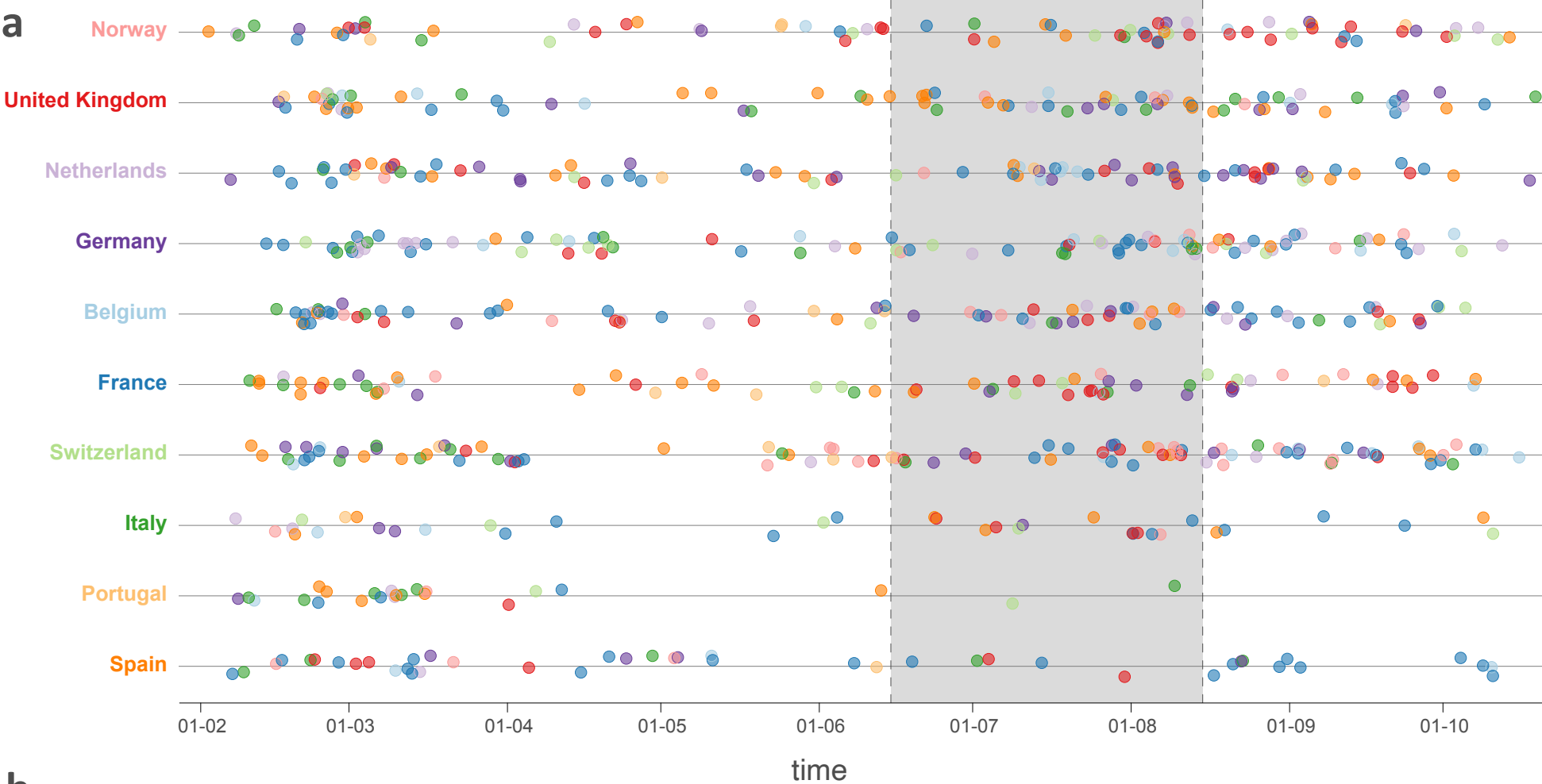
777

778 **Extended Data Table 3.** Mobility to or from each country within our 10-country sample as the percentage of the  
779 total between-country mobility for these countries within Europe.

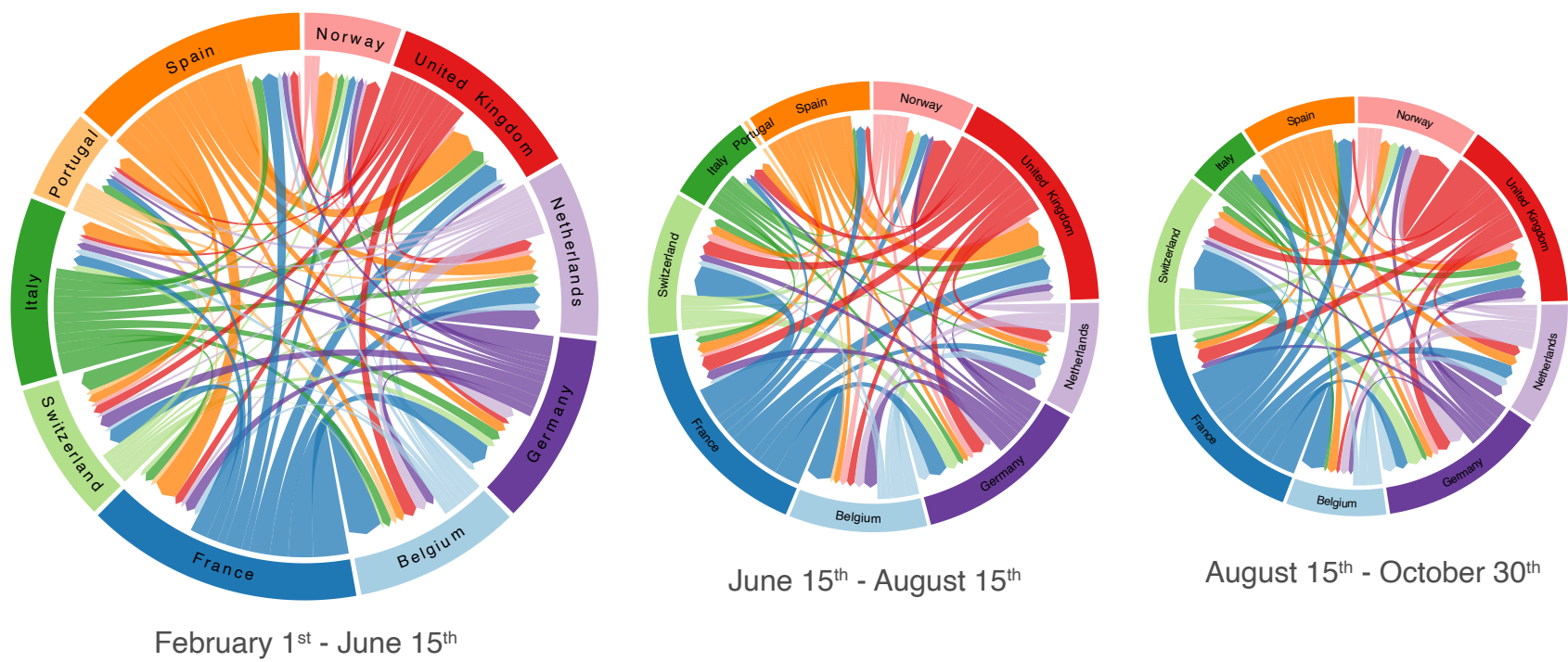
780

781 The pairwise mobility measures summarized in this table are shown in Extended Data Figure 3.

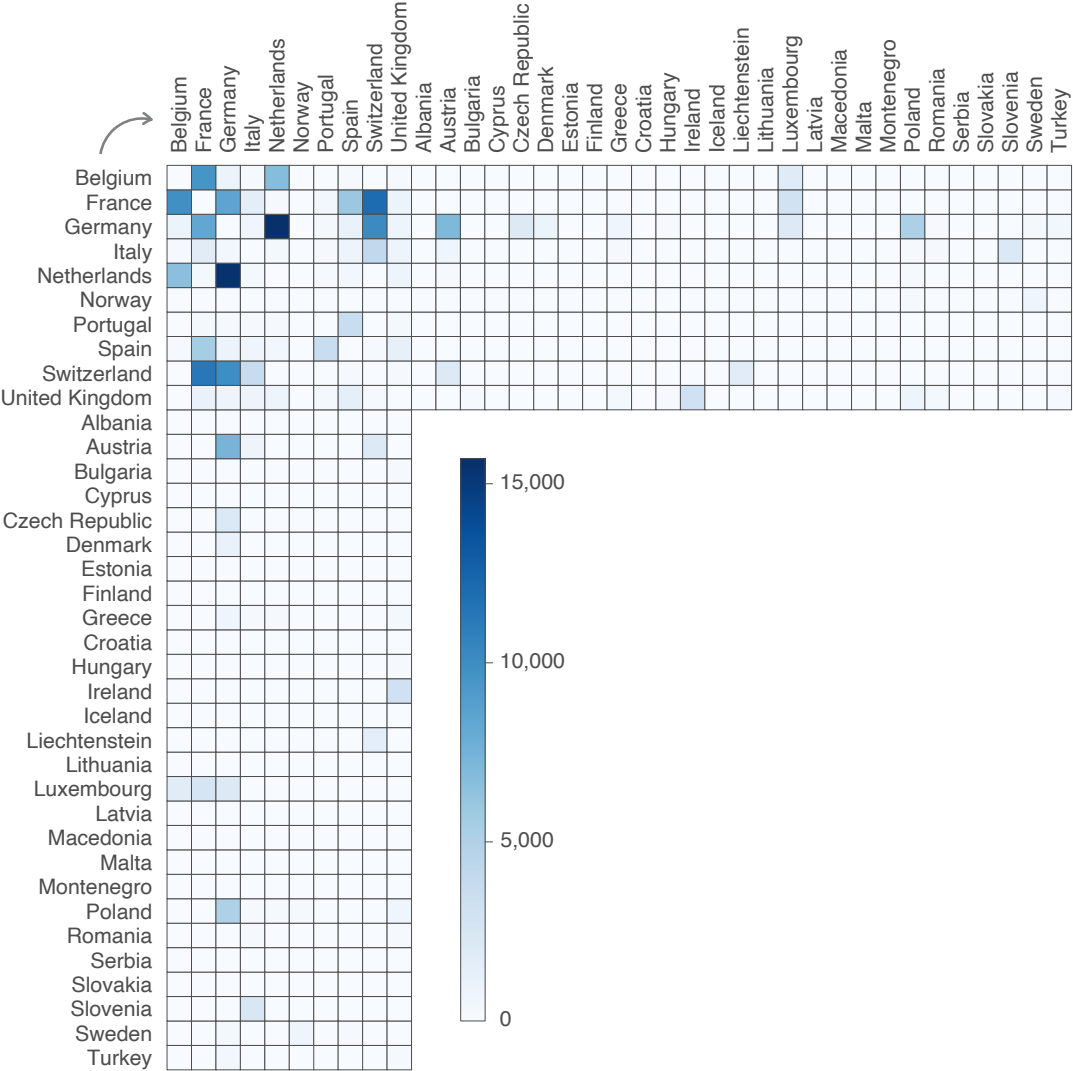
**(a) International air traffic data****(b) International Facebook mobility data (and SCI)****(c) International Google mobility data**



**b**



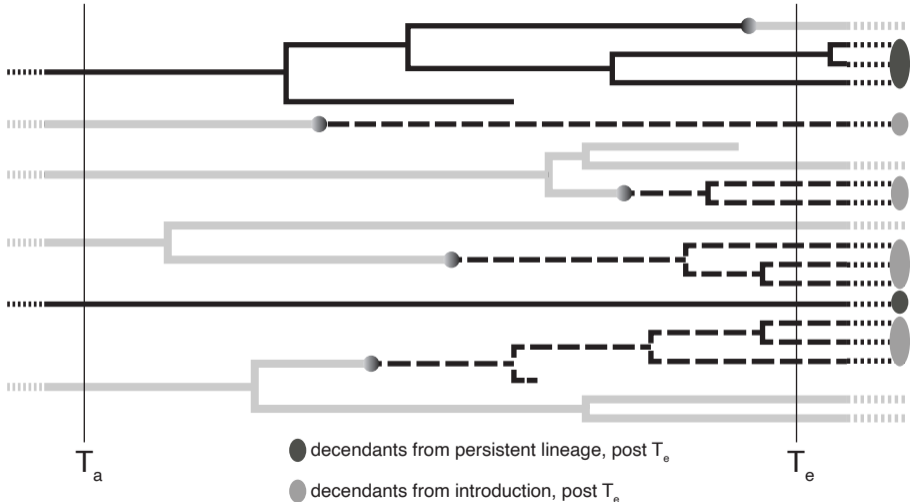




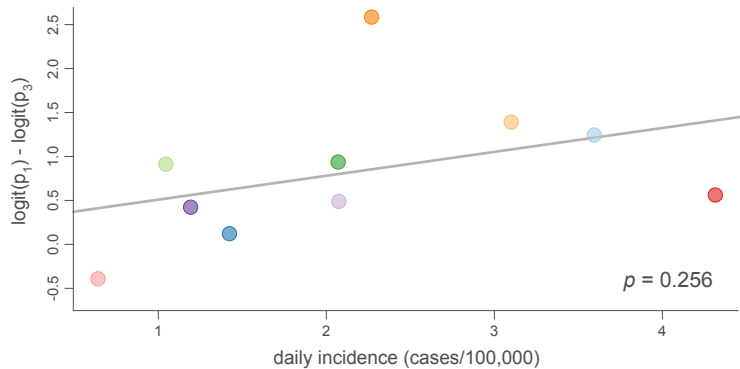
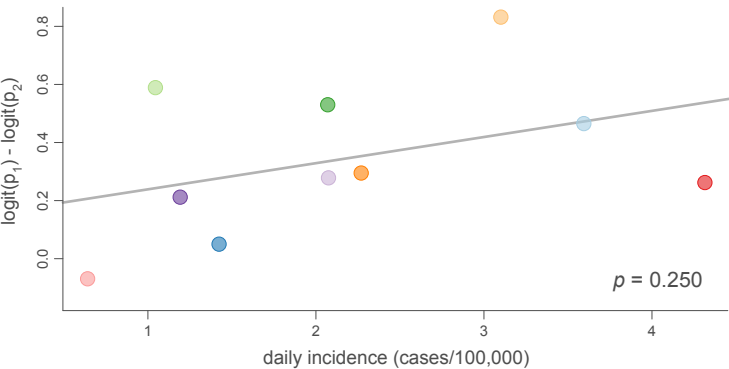
other location

persistent lineage

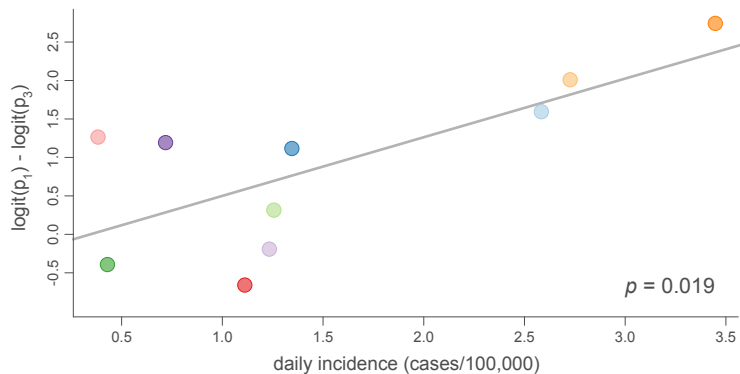
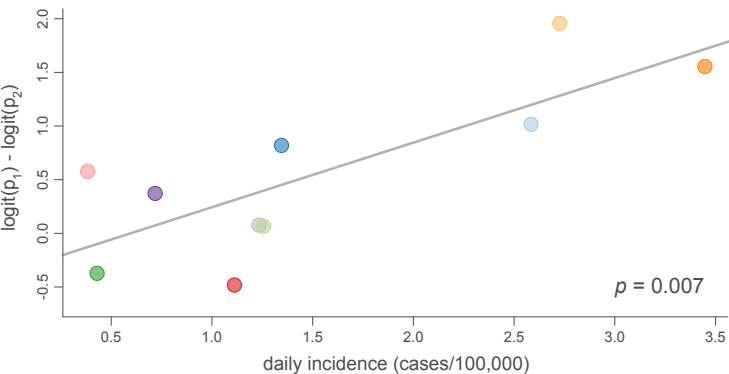
introduction



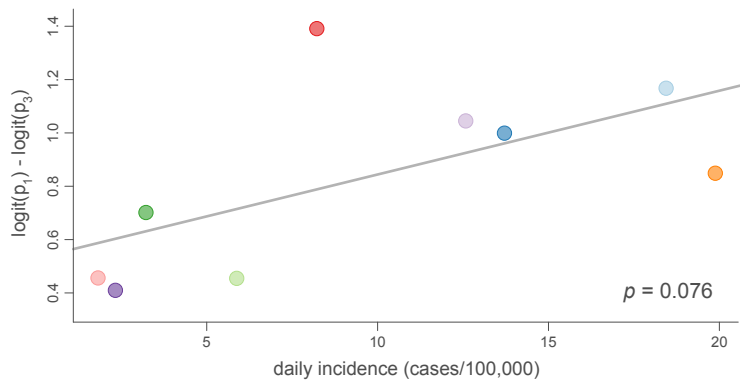
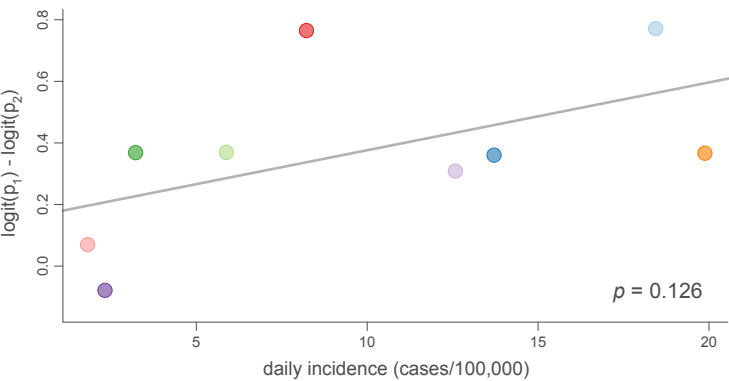
April 15<sup>th</sup> - June 15<sup>th</sup>



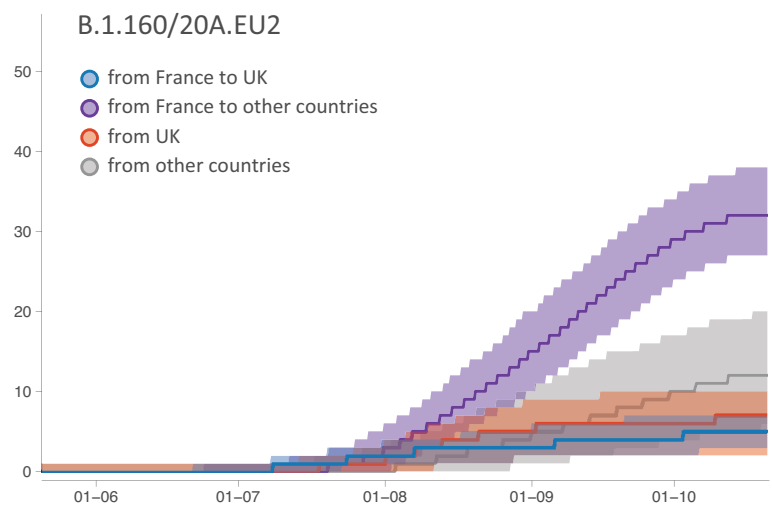
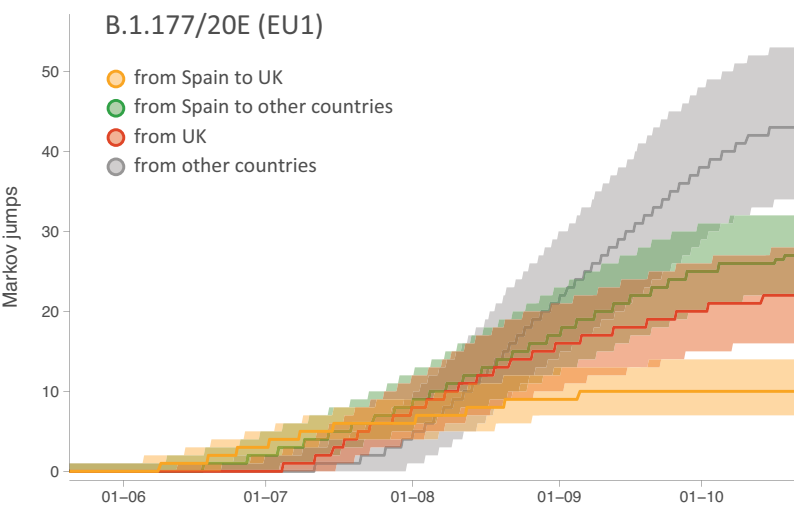
June 15<sup>th</sup> - August 15<sup>th</sup>



August 15<sup>th</sup> - October 15<sup>th</sup>







<b>country</b>	<b>genomes (Nov. 3rd, 2020)</b>	<b>genomes (Jan 5th, 2021)</b>	<b>total</b>
Belgium	183 (1091)	53 (957)	236
France	600 (1441)	167 (762)	767
Germany	246 (486)	75 (482)	321
Italy	281 (795)	75 (257)	356
The Netherlands	159 (2387)	47 (1032)	206
Norway	100 (414)	92 (482)	192
Portugal	100 (1370)	100*	200
Spain	647 (2443)	191 (827)	838
Switzerland	100 (3,019)	98 (1421)	198
The United Kingdom	493 (26,366)	152 (50,175)	645
total	2909	1050	3959

Model		Parameter estimates
Time-homogenous spatial diffusion	coalescent GLM	$\alpha = 2.604 [2.487, 2.735]$ , $\beta = 1.711 [1.603, 1.898]$
	spatial GLM	air travel: $E[\delta] = 0.018$ , $\beta( \delta=1) = 0.044 [0.001, 0.123]$ SCI: $E[\delta] = 0.004$ , $\beta( \delta=1) = 0.013 [0.003, 0.032]$ mobility: $E[\delta] > 0.999$ , $\beta( \delta=1) = 0.358 [0.258, 0.456]$
Time-inhomogeneous spatial diffusion	spatial GLM, constant inclusion probabilities	air travel: $E[\delta] = 0.018$ , $\beta( \delta=1) = 0.029 [0.001, 0.105]$ SCI: $E[\delta] = 0.008$ , $\beta \delta=1 = 0.024 [0.001, 0.078]$ mobility: $E[\delta] > 0.999$ , $\beta( \delta=1) = 0.333 [0.229, 0.438]$
	spatial GLM, time-variable inclusion probabilities	air travel: $E[\delta_i] = 0.010$ , $\beta( \delta_i=1) = 0.047 [0.002, 0.139]$ SCI: $E[\delta_i] = 0.012$ , $\beta \delta_i=1 = 0.018 [0.000, 0.062]$ mobility: $E[\delta_i] = 0.949$ , $\beta( \delta_i=1) = 0.357 [0.230, 0.503]$
	spatial GLM	mobility: $\beta = 0.271 [0.118, 0.444]$
	time-variable rate scalar GLM	mobility: $\alpha = 0.740 [0.618, 0.856]$ , $\beta = 0.504 [0.350, 0.646]$

<b>country</b>	<b>Mobility percentage</b>
Belgium	87.2
France	89.5
Germany	63.9
Italy	64.8
The Netherlands	93.2
Norway	27.1
Portugal	94.0
Spain	90.3
Switzerland	84.8
The United Kingdom	48.6