

Untangling introductions and persistence in COVID-19 resurgence in Europe

Philippe Lemey^{1,2}, Nick Ruktanonchai^{3,4}, Samuel L. Hong¹, Vittoria Colizza⁵, Chiara Poletto⁵, Frederik Van den Broeck^{1,6}, Mandev S. Gill¹, Xiang Ji⁷, Anthony Levasseur⁸, Bas B. Oude Munnink⁹, Marion Koopmans⁹, Adam Sadilek¹⁰, Shengjie Lai³, Andrew J. Tatem³, Guy Baele¹, Marc A. Suchard^{11,12,13}, Simon Dellicour^{1,14}.

¹Department of Microbiology, Immunology and Transplantation, Rega Institute, KU Leuven, Leuven, Belgium.

²Global Virus Network (GVN), Baltimore, MD, USA.

³WorldPop, School of Geography and Environmental Science, University of Southampton, Southampton SO17 1BJ, UK.

⁴Population Health Sciences, Virginia Tech, Blacksburg, VA, USA.

⁵INSERM, Sorbonne Université, Institut Pierre Louis d'Epidémiologie et de Santé Publique IPLESP, F75012 Paris, France.

⁶Department of Biomedical Sciences, Institute of Tropical Medicine, Antwerp, Belgium.

⁷Department of Mathematics, School of Science & Engineering, Tulane University, New Orleans, LA, USA

⁸Microbes, Evolution, Phylogeny and Infection, Aix-Marseille Université and Marseille Institut Universitaire de France, Marseille, France.

⁹Erasmus MC, Department of Viroscience, WHO collaborating centre for arbovirus and viral hemorrhagic fever Reference and Research, Rotterdam, the Netherlands.

¹⁰Google, Mountain View, CA, USA.

¹¹Department of Biomathematics, David Geffen School of Medicine, University of California Los Angeles, Los Angeles, CA 90095, USA.

¹²Department of Biostatistics, Fielding School of Public Health, University of California Los Angeles, Los Angeles, CA 90095, USA.

¹³Department of Human Genetics, David Geffen School of Medicine, University of California Los Angeles, Los Angeles, CA 90095, USA.

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

32 **Summary paragraph**

33 Following the first wave of SARS-CoV-2 infections in spring 2020, Europe experienced a resurgence of
34 the virus starting late summer that was deadlier and more difficult to contain. Relaxed intervention
35 measures and summer travel have been implicated as drivers of the second wave. 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 15th. The success in onward transmission of newly introduced lineages is predicted by COVID-19
41 incidence during this period. The pervasive spread of variants in summer 2020 highlights the threat of
42 viral dissemination when restrictions are lifted, and this needs to be carefully considered by strategies to
43 control the current spread of variants that are more transmissible and/or evade immunity. Our findings
44 indicate that more effective and coordinated measures are required to contain spread through cross-
45 border travel even as vaccination begins to reduce 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
49 resurgence in COVID-19 cases in late summer. The number of COVID-19 infections increased rapidly, and
50 by the end of October, it was clear that the continent was deep into a second epidemic wave. This
51 forced governments to reimpose lockdowns and social restrictions in an effort to contain the second
52 wave. While these measures have reduced infection rates across Europe ¹, several countries have
53 witnessed a stabilization at high levels or even a new surge in infections. The spread of more
54 transmissible variants, in particular B.1.1.7 (Variant of Concern 202012/01 or 20I/501Y.V1 ²), which was
55 first identified in the United Kingdom (UK), has considerably exacerbated the challenge to contain
56 COVID-19.

57

58 Already early on in the pandemic, modelling studies warned about new waves due to partial relaxation
59 of restrictions ³ or seasonal variations ⁴. By mid-April, the European Commission constructed a roadmap
60 to lifting coronavirus containment measures ⁵, recommending a cautious and coordinated manner to
61 revive social and economic activities. However, the early start of the devastating second wave
62 demonstrated that there was insufficient adherence to these measured recommendations. Cross-border
63 travel, and mass tourism in particular, has been implicated as a major instigator of the second wave.
64 Genomic surveillance demonstrated that a new variant (lineage B.1.177 ⁶, 20A.EU1 [nextstrain.org]),
65 which emerged in Spain in early summer, has spread to multiple locations in Europe ⁷. While this variant
66 quickly grew into the dominant circulating SARS-CoV-2 strain in several countries, it did not appear to be
67 associated with a higher intrinsic transmissibility ⁷.

68

69 Although it appears clear that travel had a significant impact on the second wave in Europe, it remains
70 challenging to assess how it may have restructured and reignited the epidemic in the different European
71 countries. Even without resuming travel, relaxing containment measures when low-level transmission is
72 ongoing risks the proliferation of locally circulating strains. Phylodynamic analyses may provide insights
73 into the relative importance of persistence versus the introduction of new lineages, but such analyses
74 are complicated for SARS-CoV-2 for different reasons. Phylogenetic reconstructions may be poorly
75 resolved due to the relatively limited SARS-CoV-2 sequence diversity⁸. This is further confounded by the
76 degree of genetic mixing that can be expected from unrestricted travel prior to the lockdowns in spring
77 2020.

78

79 *Mobility data predicts SARS-CoV-2 spread*

80 We analyzed SARS-CoV-2 B.1 (20A) genomes from 10 European countries for which a minimal number of
81 genomes from the second wave were already available on November 3rd, 2020. Using a two-step
82 procedure that relied on subsampling relative to country-specific case counts (cfr. Methods), we
83 compiled a data set of close to 4,000 genomes sampled between January 29th and October 30th, 2020
84 (Extended Data Table 1). In order to achieve maximum resolution in our evolutionary reconstructions,
85 we constructed a Bayesian time-measured phylogeographic model that integrates mobility and
86 epidemiological data. Our approach simultaneously infers phylogenetic history and ancestral movement
87 throughout this history while also identifying the drivers of spatial spread ⁹. We used the latter
88 functionality to determine the most appropriate mobility or connectivity measure. Specifically, we
89 considered international air transportation data, the Google COVID-19 Aggregated Mobility Research
90 Dataset (also referred to here as ‘mobility data’ for short), as well as Facebook's Social Connectedness
91 Index (SCI), as covariates of phylogeographic spread (Extended Data Figure 1). The Google mobility data

92 contains anonymized mobility flows aggregated over users who have turned on the Location History
93 setting, which is off by default (cfr. Methods). The Social Connectedness Index reflects the structure of
94 social networks and has been suggested to correlate with the geographic spread of COVID-19¹⁰. To help
95 inform the phylogenetic coalescent time distribution, we parameterized the viral population size
96 trajectories through time as a function of epidemiological case count data for the countries under
97 investigation.

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

107
108 *Dynamic cross-country transmission through time*

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

122
123 According to the proportion of introductions, we estimate more viral import than export events for
124 Norway, the Netherlands, Belgium and Switzerland throughout most of the time period under
125 investigation. According to the estimated phylogenetic entropy, these countries also experienced many
126 independent transmission chains since the epidemic started to unfold. This is consistent with country-
127 specific studies; for the first wave in Belgium for example, about 331 individual introductions were
128 estimated in the ancestry of a limited sample of 740 genomes¹¹. For Portugal, we also estimate higher
129 proportions of introductions early in the first wave but with a subsequent decline to predominantly
130 export events. France, Italy and Spain on the other hand are characterized by a relatively high viral
131 export during the first wave. The proportion of introductions remains relatively low for Italy and Spain
132 following the first wave, while in France these proportions are high from mid-June until the end of July.
133 The absolute number of transitions in our sample are however low during this time period. These
134 countries also have comparatively lower entropy values early in the epidemic, with an increase for
135 France by the start of summer and a more gradual increase over time for Italy. In Spain however, the

136 genetic complexity of SARS-CoV-2 transmission chains remains limited. In the UK and Germany, the viral
137 flow in and out of the country is initially relatively balanced. A recent large-scale genomic analysis in the
138 UK indicates that this can imply very high absolute numbers of cross-country transmissions, as more
139 than 2,800 independent introduction events were identified from the analysis of 26,181 genomes ¹².
140 Although our sample is limited compared to this analysis, our reconstructions also recover major influx
141 from Spain, France and Italy during the first wave in the UK (Extended Data Figure 2). We estimate an
142 increase in the proportion of introductions for the UK from mid-June, indicating an important viral
143 import relative to export around this time. The phylogenetic entropy also peaks around this time. In
144 Germany, the proportions increase somewhat later in summer with a concomitant rise in phylogenetic
145 entropy.

146

147 *Incidence determines the success of new introductions*

148 To assess the impact of summer travel on the second wave in the different countries, we use our
149 genomic-mobility reconstruction to estimate both the number of lineages persisting in each country and
150 the number of newly introduced lineages, and how these proliferated early in the second wave. We
151 focus on a two-month time period between June 15th, on which many EU and Schengen-area countries
152 opened their borders to other countries, and August 15th, before which the majority of holiday return
153 travel is expected for many countries. We identify the number of lineages circulating in each country on
154 August 15th, and determine whether they result from a lineage that persisted since June 15th or from a
155 unique introduction after this date (independent of the number of descendants for this lineage on
156 August 15th, Extended Data Figure 4). In Figure 2, we plot i) the ratio of these unique introductions over
157 the total unique lineages (unique introductions and persisting lineages) (p_1), ii) the proportion of
158 descendant lineages on August 15th that resulted from the unique introductions over the total
159 descendants circulating on this date (p_2), and iii) the proportion of descendant tips (sampled genomes)
160 after August 15th that resulted from the unique introductions over the total number of descendant tips
161 (p_3 , cfr. Methods and Extended Data Figure 4). We estimate a posterior mean proportion of unique
162 introductions that is close to or higher than 0.5 except for Spain and Portugal. This indicates that by
163 August 15th a relatively large fraction of circulating lineages in each country was spawned by new
164 introductions over summer. Because the B.1.177/20A.EU1 variant that was predominantly disseminated
165 through summer travel does not appear to be more transmissible ⁷, this was unlikely due to intrinsic
166 advantages of the newly introduced viruses.

167

168 The two proportions of descendants from these introductions on August 15th (p_2) and after this date (p_3)
169 measure the relative success of newly introduced lineages compared to persisting lineages, indicating
170 considerable variation in onward transmission. The country estimates are ordered according to
171 decreasing average incidence during the June 15 - August 15 time period, suggesting that incidence may
172 shape the outcome of the introductions. In countries that experienced relatively high summer incidence
173 (e.g. Spain, Portugal, Belgium and France), the introductions lead to comparatively fewer descendants
174 on August 15th or after. We find a significant overall association between incidence and the difference in
175 the logit-scaled proportion of unique introductions and the logit-scaled proportion of their descendants
176 on August 15th ($p = 0.006$) as well as between incidence and the difference in the logit-scaled proportion
177 of unique introductions and the logit-scaled proportion of descendant tips after August 15th ($p = 0.019$)
178 (Figure 2). With comparatively few descendants from introductions (Figure 2), Norway may to some
179 extent be an outlier because lineages estimated as persisting in this country could in fact be

180 introductions from other Scandinavian countries that are not represented in our genome sample. We
181 recover qualitatively similar, but more variable and statistically unsupported associations between the
182 success of introductions and incidence for the two-month time periods before and after the June 15 -
183 August 15 time period (Extended Data Figure 5). This indicates that the comparatively higher proportion
184 of introductions as well as the more stable and lower incidence between June 15th and August 15th
185 provided the ideal conditions for a process of genetic drift by which introductions were able to fuel
186 transmission.

187

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

205 Discussion

206 Our Bayesian phylogeographic approach builds on a rich history of identifying drivers of spatial spread,
207 with applications to various pathogens at different spatial scales, ranging from air transportation for
208 influenza at a global scale⁹ to gravity model transmission for Ebola in West Africa¹³. Such studies use a
209 relatively limited genomic sample to gain insights into viral transmission dynamics. This is also the case
210 in our application to SARS-CoV-2 in Europe for which we further extend the phylodynamic data
211 integration approach to confront the lack of resolution offered by SARS-CoV-2 genomic data. A
212 concerted effort in containing international spread further sets apart the COVID-19 pandemic from
213 these earlier events. For this reason, we have now incorporated variation in mobility over time to
214 account for the impact of these measures. Our reconstructions show that the composition of lineages
215 circulating towards the end of the summer was to a significant extent shaped by introductions in most of
216 the European countries. The relative success of onward transmission of the introduced lineages appears
217 to be shaped by the average summer COVID-19 incidence.

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

227
228 The pandemic exit strategy offered by vaccination programs is a source of optimism that also sparked
229 proposals by EU member states to issue vaccine passports in a bid to revive travel and rekindle the
230 economy. In addition to implementation challenges and issues of fairness, there are risks associated
231 with such strategies when immunization is incomplete, as likely will be the case for the European
232 population this summer. A recent modelling study for the United Kingdom suggests that vaccination in
233 adults alone is unlikely to completely halt the spread of COVID-19 cases and that lifting containment
234 measures early and suddenly can lead to a large wave of infections¹⁴. A gradual release of restrictions
235 was shown to be critical for minimizing the infection burden¹⁴. We believe that travel policies may be a
236 key consideration in this respect because similar conditions may arise as the ones we demonstrated to
237 provide fertile ground for viral dissemination and resurgence in 2020. This may now also involve the
238 spread of variants that evade immune responses triggered by vaccines and previous infections. Well-
239 coordinated European strategies will therefore be required to manage the spread of SARS-CoV-2 and
240 reduce future waves of infection, with hopefully a more unified implementation than hitherto observed.

241 **References**

- 242 1. COVID-19 situation update for the EU/EEA, as of week 3, updated 28 January 2021.
243 <https://www.ecdc.europa.eu/en/cases-2019-ncov-eueea>.
- 244 2. Andrew Rambaut, Nick Loman, Oliver Pybus, Wendy Barclay, Jeff Barrett, Alesandro Carabelli, Tom
245 Connor, Tom Peacock, David L Robertson, Erik Volz, on behalf of COVID-19 Genomics Consortium
246 UK (CoG-UK). Preliminary genomic characterisation of an emergent SARS-CoV-2 lineage in the UK
247 defined by a novel set of spike mutations. *virological.org* [https://virological.org/t/preliminary-](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)
248 [genomic-characterisation-of-an-emergent-sars-cov-2-lineage-in-the-uk-defined-by-a-novel-set-of-](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)
249 [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).
- 250 3. Di Domenico, L., Pullano, G., Sabbatini, C. E., Boëlle, P.-Y. & Colizza, V. Impact of lockdown on
251 COVID-19 epidemic in Île-de-France and possible exit strategies. *BMC Med.* **18**, 1–13 (2020).
- 252 4. Neher, R. A., Dyrdak, R., Druelle, V., Hodcroft, E. B. & Albert, J. Potential impact of seasonal forcing
253 on a SARS-CoV-2 pandemic. *Swiss Med. Wkly* **150**, w20224 (2020).
- 254 5. McKee, M. A European roadmap out of the covid-19 pandemic. *BMJ* **369**, m1556 (2020).
- 255 6. Rambaut, A. *et al.* A dynamic nomenclature proposal for SARS-CoV-2 lineages to assist genomic
256 epidemiology. *Nat Microbiol* **5**, 1403–1407 (2020).
- 257 7. Hodcroft, E. B. *et al.* Emergence and spread of a SARS-CoV-2 variant through Europe in the summer
258 of 2020. *medRxiv* (2020) doi:10.1101/2020.10.25.20219063.
- 259 8. Morel, B. *et al.* Phylogenetic analysis of SARS-CoV-2 data is difficult. *Molecular Biology and*
260 *Evolution* (2020) doi:10.1093/molbev/msaa314.
- 261 9. Lemey, P. *et al.* Unifying Viral Genetics and Human Transportation Data to Predict the Global
262 Transmission Dynamics of Human Influenza H3N2. *PLoS Pathog.* **10**, e1003932 (2014).
- 263 10. Kuchler, T., Russel, D. & Stroebel, J. The Geographic Spread of COVID-19 Correlates with the
264 Structure of Social Networks as Measured by Facebook. (2020) doi:10.3386/w26990.
- 265 11. Dellicour, S. *et al.* A Phylodynamic Workflow to Rapidly Gain Insights into the Dispersal History and
266 Dynamics of SARS-CoV-2 Lineages. *Mol. Biol. Evol.* (2020) doi:10.1093/molbev/msaa284.
- 267 12. du Plessis, L. *et al.* Establishment and lineage dynamics of the SARS-CoV-2 epidemic in the UK.
268 *Science* (2021) doi:10.1126/science.abf2946.
- 269 13. Dudas, G. *et al.* Virus genomes reveal factors that spread and sustained the Ebola epidemic. *Nature*
270 **544**, 309–315 (2017).
- 271 14. Moore, S., Hill, E. M., Tildesley, M. J., Dyson, L. & Keeling, M. J. Vaccination and non-pharmaceutical
272 interventions for COVID-19: a mathematical modelling study. *Lancet Infect. Dis.* (2021)
273 doi:10.1016/S1473-3099(21)00143-2.

274 Figure Legends

275

276 **Figure 1. Mobility, genome sampling, case counts and phylogeographic summaries through time for 10 European countries.**

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

290

291 **Figure 2. Posterior estimates for the relative importance of lineage introduction events in 10 European countries and their association with incidence. (a)**

292 We report three summaries (posterior mean and 95% HPD intervals) for each country: the ratio
293 of unique introductions over the total number of unique persisting lineages and unique introductions between June 15 and
294 August 15, 2020 (p_1), the ratio of descendant lineages from these unique introduction events over the total number of
295 descendants circulating on August 15, 2020 (p_2), and the ratio of descendant taxa from these unique introductions over the
296 total number of descendant taxa sampled after August 15, 2020 (p_3)(cfr. Extended Data Figure 4). The dot sizes are proportional
297 to: (1) the total number of unique lineage introductions identified between June 15 and August 15, 2020, (2) the total number
298 of lineages inferred on August 15, 2020, and (3) the total number of descendant sequences after August 15, 2020. The third
299 ratio is not included for Portugal due to insufficient sequences sampled after August 15, 2020. The lower panels plot the
300 difference between $\text{logit}(p_1)$ and $\text{logit}(p_2)$ against incidence between June 15th and August 15th (b), and the difference between
301 $\text{logit}(p_1)$ and $\text{logit}(p_3)$ against incidence for the same period (c). A large difference in logit-scaled proportions represents a strong
302 decline in these proportions.

303

304 **Figure 3.** Phylogeographic estimates of SARS-CoV-2 spread in 10 European countries. The radial tree in the center represents
305 the maximum clade credibility tree summary of the Bayesian inference. Colors correspond to the countries in the legend. The
306 two clades corresponding to B1.1777/20A.EU1 and B1.160/20A.EU2 are highlighted in grey. The circular migration flow plots for
307 these variants are based on the posterior expectations of the Markov jumps. In these plots, migration flow out of a particular
308 location starts close to the outer ring and ends with an arrowhead more distant from the destination location. For
309 B1.1777/20A.EU1, we summarize phylogeographic transitions as posterior mean estimates with 95% HPD intervals over time
310 for four types of Markov jumps: i) from Spain to the UK, ii) from Spain to other countries, iii) from the UK, and iv) from other
311 countries.

312

313 Extended Data Figures

314

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

318

319 **Extended Data Figure 2.** Estimated introductions through time in the 10 European countries and circular migration
320 flow plots summarizing the estimated transitions between the countries for different time intervals throughout the
321 SARS-CoV-2 evolutionary history. (a) The introductions through time serve as an illustration and are based on the
322 Markov jump history in the MCC tree. We note that the posterior distribution of trees is accompanied with
323 considerable uncertainty about the location of origin, destination and timing of the transitions, which is difficult to
324 appropriately visualize. The grey box represents the time period from June 15th to August 15th. (b) The circular
325 migration flow plots are based on the posterior expectations of the Markov jumps. The size of the plots reflects the

326 total number of transitions for each period. In these plots, migration flow out of a particular location starts close to
327 the outer ring and ends with an arrowhead more distant from the destination location.

328

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

332

333 **Extended Data Figure 4.** Conceptual representation of persistent lineages and introductions during the time
334 interval delineated by the evaluation time (T_e) and the ancestral time (T_a). At T_e , we evaluate how many lineages
335 are circulating in the location of interest, in this case 12 (lineages in other locations are represented by thick grey
336 branches). We subsequently identify whether these lineages maintained this location up to T_a in their ancestry or
337 whether they result from an introduction event in the time interval of interest. By determining whether other
338 lineages circulating in the location of interest at T_e are descendants of the same persistent lineage or whether they
339 share an introduction event, we identify the unique persistent lineages or introductions, in this case 2 and 4
340 respectively. In addition to the proportion of unique introductions (4/6), we also summarize the proportion of their
341 descendants at T_e (9/(9+3) in this case) and the proportion of their descendants in terms of sampled tips after T_e .
342 Those tips are not shown here but conceptually represented for both introductions and persistent lineages by
343 ovals.

344

345 **Extended Data Figure 5.** Scatter plots of the difference in the logit proportion of unique introductions (p_1) and the
346 logit proportion of their descendants on August 15th (p_2) against incidence and the difference in the logit
347 proportion of unique introductions and the logit proportion of descendant tips after August 15th (p_3) against
348 incidence. Both plots are shown for the period between April 15th and June 15th and for the period between August
349 15th and October 15th respectively. Linear regression p -values are included in the lower right corner of the plots.

350

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

354

355 **Extended Data Figure 7.** Phylogeographic transitions for lineages B1.1777/20A.EU1 and B1.160/20A.EU2.
356 Cumulative phylogeographic transitions are summarized as posterior mean estimates with 95% HPD intervals over
357 time for 4 types of Markov jumps. For B1.1777/20A.EU1: i) from Spain to the UK, ii) from Spain to other countries,
358 iii) from the UK, and iv) from other countries; For B1.160/20A.EU2: i) from France to the UK, ii) from France to
359 other countries, iii) from the UK, and iv) from other countries.

360

361 **Extended Data Figure 8.** Posterior summary of the GLM random effects. The posterior distributions for all random
362 effects in log space are summarized as an error bar plot. The mean effect sizes are represented by white horizontal
363 lines, the boxes show interquartile ranges and the whiskers represent the full ranges of the posterior distributions.

364

365 **Extended Data Figure 9.** Comparison between Google and Facebook aggregate international mobility data. We
366 summarize monthly correlations using scatter plots and Spearman's rank correlation. Each dot in the scatter plots
367 corresponds to a specific pair of European countries considered in our study.

368

369 **Extended Data Figure 10.** Root-to-tip divergence as a function of sampling time for the 3,959 genome data set
370 with a different rooting of the same maximum likelihood tree. A. Tree rooted according to the best-fitting root
371 under the heuristic residual mean squared criterion. B. Tree rooted along the branch leading to the cluster of three
372 Bavarian genomes that resulted from an independent introduction into Europe.

373

374 **Extended Data Tables**

375 **Extended Data Table 1.** Genome sampling by country, collected on Nov. 3rd, 2020, and updated on Jan 5th, 2021.

376

377 **Extended Data Table 2.** Parameter estimates for the various Bayesian time-measured phylogeographic models
378 applied to the 3,959 genome data set.

379

380 **Extended Data Table 3.** Mobility to or from each country within our 10-country sample as the percentage of the
381 total between-country mobility within Europe.

382

383 **Extended Data Table 4.** Log marginal likelihood estimates for case count covariates with different lags in the
384 coalescent GLM model.

385

386 **Extended Data Table 5.** Parameter estimates for phylogeographic GLM analyses including an origin and destination
387 covariate based on the residuals for a regression analysis between the number of genomes and case counts.

388

389 **Methods**

390 *Sequence data and subsampling*

391 We used a two-step genome data collection procedure. We first evaluated the available genomes from
392 European countries in GISAID ¹⁵ on November 3, 2020. We selected genomes from Belgium, France,
393 Germany, Italy, Netherlands, Norway, Portugal, Spain, Switzerland and the UK primarily based on the
394 availability of genome data from both the first and second wave at that time but also because of their
395 high ratio of genomes to positive cases. A total of 39,812 genomes were available for these countries on
396 November 3, 2020; the available number of genomes by country are listed in Extended Data Table 1.
397 Portugal represented an exception because data for this country were limited to the first wave at that
398 time, but we included genomes from Portugal because of its potential importance as a summer travel
399 location.

400
401 We aligned the genomes from each country using MAFFT v7.453 ¹⁶ and trimmed the 5' and 3' ends and
402 only retained unique sequences from each location. To further mitigate the disparities in sampling, we
403 subsampled each country proportionally to the cumulative number of cases on October 21st (the most
404 recently sampled sequence at the time) by setting an arbitrary threshold of 6.5 sequences per 10,000
405 cases, with a minimum number of 100 sequences per country. To maximize the temporal and spatial
406 coverage in each country, we binned genomes by epi-week and sampled as evenly as possible, sampling
407 from a different region within the country when available. Only sequences from the B.1 lineage with the
408 D614G mutation and exact sampling dates were selected for the analyses. From the final aligned
409 sequence set, we removed 12 potential outliers, based on a root-to-tip regression on TempEst v1.5.3 ¹⁷
410 on a maximum-likelihood tree inferred with IQTREE v2.0.3 ¹⁸, yielding a data set of 2,909 genomes
411 (Extended Data Table 1).

412
413 Because of the nature of genome sequence accumulation, fewer recently sampled genomes were
414 available for most countries on November 3rd (relative to the case counts at this time). Because our
415 primary goal was to assess the persistence and introduction of lineages leading up to the second wave,
416 we sought to augment our data set with more recent genomes, having already performed analyses on
417 the initial data set. In the section on Bayesian evolutionary reconstructions, we outline how we update
418 these analyses accordingly. On January 5th, 2021, we updated our dataset by adding over 1,000 non-
419 identical sequences collected between August 1st and October 31st (out of a total of 56,395 available
420 genomes; the available and selected number of genomes by country are listed in Extended Data Table
421 1). For Portugal, we extended this period back to June 22nd (the most recent sampling date for the
422 previous Portuguese selection). We downloaded all new B.1 sequences with the D614G mutation
423 collected during the selected time period from GISAID and performed the following subsampling. The
424 number of genomes to add by country was obtained by raising the threshold ratio of sequences/cases to
425 8.5 and increasing the minimum number of sequences to 200. To bias the temporal coverage towards
426 more recent samples, the genomes from each country were binned by week and sampled such that the
427 number of sequences added by week was proportional to an exponential function of the form $e^{t/4}$,
428 where $t=0$ represents August 1st and $t=13$ is October 31st. For Portugal, we did not use this preferential
429 sampling as we needed to include close to all available genomes to raise the number of genomes to 200.
430 The sampled sequences were then deduplicated and outliers were removed as described in the previous

431 section. With the additional selection of 1,050 genomes, we arrived at a data set of 3,959 genomes
432 (Extended Data Table 1).

433

434 *Mobility data*

435 We analysed four different mobility/connectivity measures: air traffic flows, a social connectedness
436 index provided by Facebook, as well as aggregate Facebook¹⁹ and Google international mobility data.
437 Air traffic flow data were obtained from the International Air Transport Association
438 (<http://www.iata.org>) and based on the number of origin-destination tickets while also taking into
439 account connections at intermediate airports²⁰. We used monthly air traffic data between the 10
440 European countries under investigation for the time period between January 2020 and October 2020.
441 The social connectedness index (SCI) is an anonymized snapshot of active Facebook users and their
442 friendship networks to measure the intensity of social connectedness between countries
443 (<https://data.humdata.org/>)²¹. In practice, the SCI measures the relative probability of a Facebook
444 friendship link between two users of the application in different countries. We used the SCI calculated
445 for the 10 European countries represented in our genomic sample as of August 2020.

446

447 The Google COVID-19 Aggregated Mobility Research Dataset contains anonymized mobility flows
448 aggregated over users who have turned on the Location History setting (on a range of platforms²²),
449 which is off by default. To produce this dataset, machine learning is applied to logs data to automatically
450 segment it into semantic trips²³. To provide strong privacy guarantees, all trips were anonymized and
451 aggregated using a differentially private mechanism²⁴ to aggregate flows over time (see
452 <https://policies.google.com/technologies/anonymization>). This research was done on the resulting
453 heavily aggregated and differentially private data. No individual user data was ever manually inspected,
454 only heavily aggregated flows of large populations were handled. All anonymized trips were processed
455 in aggregate to extract their origin and destination location and time. For example, if users traveled from
456 location a to location b within time interval t , the corresponding cell (a, b, t) in the tensor would be $n \pm$
457 η , where η is Laplacian noise. The automated Laplace mechanism adds random noise drawn from a zero-
458 mean Laplace distribution and yields (ϵ, δ) -differential privacy guarantee of $\epsilon = 0.66$ and $\delta = 2.1 \times 10^{-29}$
459 per metric. The parameter ϵ controls the noise intensity in terms of its variance, while δ represents the
460 deviation from pure ϵ -privacy. The closer they are to zero, the stronger the privacy guarantees. We used
461 aggregated mobility flows between the 10 European countries and summarized them by two-week or
462 monthly time periods between January 2020 and October 2020.

463

464 Finally, we also considered international mobility data from Facebook mobility data as an alternative to
465 Google mobility data. These data are based on numbers of Facebook users moving over large distances,
466 like air or train travel. Counts of international travel patterns are updated daily based only on users who
467 have opted to share precise location data from their device with the Facebook mobile app through
468 location services. Also in this case, we used aggregated mobility flows between the 10 European
469 countries and summarized them by month between January 2020 and October 2020. Because
470 international aggregate mobility data obtained from Google and Facebook are highly correlated
471 (monthly Spearman correlation ranging from 0.84 to 0.92; Extended Figure 9), we only included the
472 Google aggregate mobility data as a covariate in the phylogeographic analyses. We note that the
473 mobility data are subject to limitations as these may not be representative for the population as whole
474 and their representativeness may vary by location.

475

476 *Bayesian evolutionary reconstructions*

477 - Joint sequence-trait inference with a time-homogeneous GLM diffusion model

478 We performed Bayesian evolutionary reconstruction of timed phylogeographic history using BEAST 1.10
479 ²⁵ incorporating genome sequences, their country and date of sampling, epidemiological and
480 mobility/connectivity data. Because of the relatively low degree of resolution offered by the sequence
481 data, our full probabilistic model specification focuses on i) relatively simple model specifications and ii)
482 informing parameters by additional non-genetic data sources. We modeled sequence evolution using an
483 HKY85 nucleotide substitution model with gamma-distributed rate variation among sites and a strict
484 molecular clock model. Our genome set includes three genomes from an early outbreak in Bavaria,
485 which was caused by an independent introduction from China ^{26,27}. We therefore constrained these
486 genomes as an outgroup in the analysis, which according to root-to-tip regression plots as a function of
487 sampling time resulted in a better correlation coefficient/R-squared compared to the best-fitting root
488 under the heuristic mean residual squared criterion (Extended Figure 10)¹⁷.

489

490 As a coalescent tree prior, we modeled the effective population size trajectory as a piecewise constant
491 function that changes values at pre-specified times (following ²⁸), with log population sizes modelled as
492 a deterministic function of log COVID-19 case counts (following ²⁹). This reduces the nonparametric
493 skygrid parameterization to a generalized linear model (GLM) formulation with an estimable regression
494 intercept (α) and coefficient (β). In this parameterization, a coefficient estimate centered around 0
495 would imply constant population size dynamics through time. We specified two-week intervals and
496 summarized as a covariate the total case counts over these time intervals for the 10 countries of
497 sampling (obtained from <https://www.ecdc.europa.eu/en/covid-19/data>). The earliest interval with
498 non-zero cases counts was from 2020-01-14 to 2020-01-28; before 2020-01-14, the log-transformed and
499 standardized case count covariate was set to the equivalent of 1 case. We also tested whether a lag-
500 time was required for the case count covariate using marginal likelihood estimation (MLE) ³⁰.
501 Specifically, we shifted the case counts by 1, 2, 3 and 4 weeks before summarizing them according to
502 two-week intervals and estimated the model fit of these covariates against case counts without lag time
503 (Supporting data Table 3). To mitigate the computational burden associated with the MLE procedure,
504 we performed these analyses on a subset of 1,000 genomes (obtained using the phylogenetic diversity
505 analyzer tool ³¹). We estimated the highest (log) marginal likelihood for a two-week lag time (Extended
506 Data Table 4) and used this for the case count covariate in our analyses.

507

508 Similar to sequence evolution, we modelled the process of transitioning through discrete location states
509 (countries of sampling) according to a continuous-time Markov chain (CTMC) ³². We employed a
510 parameterization that models the log transition rates as a log linear function of mobility/connectivity
511 covariates ⁹. The Bayesian implementation of this model simultaneously estimates phylogenetic history,
512 ancestral movement and the contribution of covariates to the movement patterns ⁹. While we mainly
513 use this approach to obtain well-informed phylodynamic estimates, we also make use of its capacity to
514 identify the most relevant mobility measure to inform our reconstructions. As covariates we considered
515 Facebook's SCI, air transportation data and mobility data. For the two time-variable mobility measures,
516 we used the average of the log-transformed and standardized monthly mobility measures as a single
517 covariate in our time-homogeneous phylogeographic GLM model. In this GLM formulation, we estimate

518 positive effect sizes for each covariate as well as their inclusion probability through a spike-and-slab
519 procedure⁹. Although we subsampled the number of SARS-CoV-2 genomes by country in proportion to
520 case counts, they do not fully correspond because we used a minimum number of genomes for
521 countries with low case counts. We therefore evaluated whether this resulted in signal for sampling bias
522 by including an origin and destination covariate in the GLM based on the residuals for a regression
523 analysis between genomes and case counts (following¹³). We performed this analysis using a set of
524 empirical trees (cfr. below) in both a time-homogeneous and time-inhomogeneous model, but found no
525 support for these additional covariates (Extended Data Table 5).

526
527 We performed inference under the full model specification using Markov chain Monte Carlo (MCMC)
528 sampling and used the BEAGLE library v3³³ to increase computational performance. We specified
529 standard transition kernels on all parameters, except for the regression coefficients of the piecewise-
530 constant coalescent GLM model. For these parameters, we implemented new Hamiltonian Monte Carlo
531 (HMC) transition kernels to improve sampling efficiency. These kernels use principles from Hamiltonian
532 dynamics and their approximate energy conserving properties to reduce correlation between successive
533 sampled states, but require computation of the gradient of the model log-posterior with respect to the
534 parameters of interest, in addition to efficient evaluation of the log-posterior that BEAGLE provides. To
535 accomplish this, we extended our previous analytic derivation of the gradient of the log-density from the
536 skygrid coalescent model with respect to the log-population-sizes³⁴ to now be with respect to the
537 regression coefficients using the chain rule and their regression design matrix.

538
539 Due to the data set size, MCMC burn-in takes up considerable computational time. We therefore
540 iterated through a series BEAST inferences, initially only considering sequence evolution and
541 subsequently adding the location data, to arrive at a tree distribution from which trees were taken as
542 starting trees in our final analyses. The latter was composed of multiple independent MCMC runs that
543 were run sufficiently long to ensure that their combined posterior samples achieved effective sample
544 sizes (ESSs) larger than 100 for all continuous parameters.

545
546 - Data augmentation through online BEAST

547 As we updated our dataset following initial analyses of the 2,909 genome collection using the approach
548 discussed in the previous subsection, we sought to capitalize on these efforts to limit the burn-in for
549 subsequent analyses of the 3,959 dataset. Specifically, we adopted the distance-based procedure to
550 insert new taxa into a time-measured phylogenetic tree sample as implemented in the BEAST
551 framework for online inference³⁵. We subsequently use the augmented tree as the starting tree for the
552 analyses of the updated dataset.

553
554 - Time-inhomogeneous reconstructions

555 To accommodate the time-variability of the mobility measures, we constructed epoch model extensions
556 of the discrete phylogeography approach that allow specifying arbitrary intervals over the evolutionary
557 history and associating them with different model parameterizations³⁶. As a complement to testing
558 covariates of spatial diffusion using a time-homogeneous model, we used the epoch extension to specify
559 monthly intervals allowing us to incorporate monthly mobility matrices (air transportation data were
560 only available as monthly numbers), but assuming time-homogeneous effect sizes and inclusion

561 probabilities. Monthly covariates were again log-transformed and standardized after adding a pseudo-
562 count to each entry in the monthly matrices.

563

564 In addition, we performed another analysis in which we relaxed the constant-through-time inclusion
565 probability of the covariates. In this model specification, each interval is associated with a specific set of
566 indicator variables to represent the inclusion/exclusion of covariates, but we pool information about
567 predictor inclusion across the intervals using hierarchical graph modelling³⁷. This approach uses a set of
568 indicator variables to model covariate inclusion at the hierarchical level but allows interval-specific
569 inclusion or predictors to diverge from the hierarchical level with a non-zero probability (with the
570 number of differences modelled as a binomial distribution,³⁷), which was set to 0.10 in our case. We
571 estimated hierarchical and interval-level inclusion using spike-and-slab³⁷.

572

573 Finally, we performed an analysis using the time-inhomogeneous model in which the interval-specific
574 transition rates are modelled as a function of the single covariate that is supported by the analyses
575 above leveraging aggregate mobility. We incorporated more variability through time by specifying two-
576 week intervals (similar to the coalescent GLM interval specification). In addition, we add time-
577 homogeneous random effects to the phylogeographic transition rate parameterization in order to
578 account for potential biases in the ability of mobility to predict phylogeographic spread. While posterior
579 mean estimates for these random effects vary, only very few indicate that individual phylogeographic
580 transition rates significantly deviate from the mobility data (Extended Data Figure 8). The time-
581 inhomogeneous GLM approach we employ allows modelling relative differences in transition rates, but
582 also the overall rate of migration between countries varies through time, and importantly, this is
583 strongly impacted by intervention strategies. To accommodate these dynamics, we further extended
584 this model by incorporating a time-inhomogeneous overall CTMC rate scaler and parameterize it as a log
585 linear function of the total monthly between-country log-transformed and standardized mobility (time-
586 variable rate scalar GLM in Extended Data Table 2). To generate realisations of the discrete location
587 CTMC process and obtain estimates of the transitions (Markov jumps) between states under this model,
588 we employed posterior inference of the complete Markov jump history through time^{9,38}.

589

590 While the epoch model allows us to flexibly accommodate time-variable spatial dynamics, it
591 considerably increases the computational burden associated with likelihood evaluations. In order to
592 efficiently draw inference under this model for our large data set, we fit the time-inhomogeneous
593 spatial diffusion process to a set of trees inferred under the time-homogeneous GLM diffusion model
594 described above. Although likelihood evaluations remain computationally expensive, even with the
595 speed-up offered by GPU computation with BEAGLE, eliminating simultaneous tree estimation
596 tremendously reduces parameter-space, requiring only modest MCMC chain lengths to adequately
597 explore it.

598

599 - Posterior Summaries

600 We assessed MCMC mixing (e.g. using ESSs) and summarized continuous parameter estimates using
601 Tracer v1.7.1³⁹. Credible intervals were computed as 95% HPD intervals. Trees were visualized using
602 FigTree v1.4.4 (available at <https://github.com/rambaut/figtree/releases>). In terms of phylogeographic
603 estimates, we mainly focused on i) transitions to each location and from each location (based on Markov
604 jump estimates) instead of pairwise transitions, ii) ratios of these transitions and iii) how these

605 transitions structured transmission chains in individual countries. Transitions to each and from each
606 location avoid drawing conclusions about direct migration between countries, which can be tenuous
607 given the incomplete genomes coverage of Europe, while their ratios avoid using absolute numbers of
608 transitions, which are highly sample-dependent. Phylogeographic inference is limited to reconstructing
609 the transitions in the ancestral history of a sample of sequences, which will only be a small fraction of
610 the actual migration events especially when these events result in insufficient onward transmission to
611 be captured in our limited sample. In addition, SARS-CoV-2 genome data can be poorly resolved and
612 identical genomes in different locations are consistent with hypotheses that involve both a sparse and a
613 rich number of virus flows between these locations. As the data hold little information to distinguish
614 these hypotheses, we only consider sparse scenario's by including only unique sequences for each
615 location. A joint inference of sequence evolution and discrete spatial diffusion would err on the side of
616 sparse hypotheses anyway because it will tend to cluster identical sequences that share a location.
617 Despite the general underestimation of spatial dispersal, a phylogeographic inference is still likely to
618 capture the transition events with important onward transmission, and evaluating the importance of
619 such events relative to persistence is a major focus of this study. Cryptic transmission also complicates
620 the ability to reconstruct spatial dispersal, but we expect this to be equally likely for introductions and
621 persistence and therefore focus on their ratio for each location.

622

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

626

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

630

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

638

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

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

641

642 • *PersistenceSummarizer* also uses posterior tree distributions annotated with Markov jumps to
643 summarize the number of lineages at a particular point in time (evaluation time, T_e , cfr.
644 Extended Figure 5), which location states they are associated with, since what time point in the
645 past they have maintained that state and how many sampled descendants they have after time
646 T_e (Extended Figure 5). In addition, it allows identifying how long these lineages have circulated
647 independently prior to T_e , so before sharing common ancestry with other lineages that
648 maintained the same location state. This information allows us to determine how many lineages

649 are circulating at T_e that stem either from a unique persistent lineage (maintaining the same
650 location states) or unique introduction event since a particular time prior to T_e . The association
651 between incidence and the difference in the logit proportion of unique introductions and the
652 logit proportion of their descendants on August 15th was evaluated using a p -value obtained by
653 a linear regression analysis.

- 654
- 655 15. Shu, Y. & McCauley, J. GISAID: Global initiative on sharing all influenza data - from vision to reality.
656 *Euro Surveill.* **22**, (2017).
 - 657 16. Katoh, K., Asimenos, G. & Toh, H. Multiple alignment of DNA sequences with MAFFT. *Methods Mol.*
658 *Biol.* **537**, 39–64 (2009).
 - 659 17. Rambaut, A., Lam, T. T., Max Carvalho, L. & Pybus, O. G. Exploring the temporal structure of
660 heterochronous sequences using TempEst (formerly Path-O-Gen). *Virus Evol* **2**, vew007 (2016).
 - 661 18. Minh, B. Q. *et al.* IQ-TREE 2: New Models and Efficient Methods for Phylogenetic Inference in the
662 Genomic Era. *Mol. Biol. Evol.* **37**, 1530–1534 (2020).
 - 663 19. Maas, P. Facebook disaster maps: Aggregate insights for crisis response & recovery. in *Proceedings*
664 *of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining (ACM,*
665 *2019)*. doi:10.1145/3292500.3340412.
 - 666 20. Gilbert, M. *et al.* Preparedness and vulnerability of African countries against importations of COVID-
667 19: a modelling study. *Lancet* **395**, 871–877 (2020).
 - 668 21. Bailey, M., Cao, R., Kuchler, T., Stroebel, J. & Wong, A. Social Connectedness: Measurement,
669 Determinants, and Effects. *J. Econ. Perspect.* **32**, 259–280 (2018).
 - 670 22. Kraemer, M. U. G. *et al.* Mapping global variation in human mobility. *Nat Hum Behav* **4**, 800–810
671 (2020).
 - 672 23. Bassolas, A. *et al.* Hierarchical organization of urban mobility and its connection with city livability.
673 *Nat. Commun.* **10**, 4817 (2019).
 - 674 24. Wilson, R. J. *et al.* Differentially Private SQL with Bounded User Contribution. *Proceedings on*
675 *Privacy Enhancing Technologies* **2020**, 230–250.
 - 676 25. Suchard, M. A. *et al.* Bayesian phylogenetic and phylodynamic data integration using BEAST 1.10.
677 *Virus Evol* **4**, vey016 (2018).
 - 678 26. Böhmer, M. M. *et al.* Investigation of a COVID-19 outbreak in Germany resulting from a single
679 travel-associated primary case: a case series. *Lancet Infect. Dis.* **20**, 920–928 (2020).
 - 680 27. Worobey, M. *et al.* The emergence of SARS-CoV-2 in Europe and North America. *Science* **370**, 564–
681 570 (2020).
 - 682 28. Gill, M. S. *et al.* Improving Bayesian population dynamics inference: a coalescent-based model for
683 multiple loci. *Mol. Biol. Evol.* **30**, 713–724 (2013).
 - 684 29. Faria, N. R. *et al.* Genomic and epidemiological monitoring of yellow fever virus transmission
685 potential. *Science* **361**, 894–899 (2018).
 - 686 30. Baele, G., Lemey, P. & Suchard, M. A. Genealogical Working Distributions for Bayesian Model
687 Testing with Phylogenetic Uncertainty. *Systematic Biology* vol. 65 250–264 (2016).
 - 688 31. Chernomor, O. *et al.* Split diversity in constrained conservation prioritization using integer linear
689 programming. *Methods Ecol. Evol.* **6**, 83–91 (2015).
 - 690 32. Lemey, P., Rambaut, A., Drummond, A. J. & Suchard, M. A. Bayesian phylogeography finds its roots.
691 *PLoS Comput. Biol.* **5**, e1000520 (2009).
 - 692 33. Ayres, D. L. *et al.* BEAGLE 3: Improved performance, scaling, and usability for a high-performance

- 693 computing library for statistical phylogenetics. *Syst. Biol.* **68**, 1052–1061 (2019).
- 694 34. Baele, G., Gill, M. S., Lemey, P. & Suchard, M. A. Hamiltonian Monte Carlo sampling to estimate
695 past population dynamics using the skygrid coalescent model in a Bayesian phylogenetics
696 framework. *Wellcome Open Res* **5**, 53 (2020).
- 697 35. Gill, M. S., Lemey, P., Suchard, M. A., Rambaut, A. & Baele, G. Online Bayesian Phylodynamic
698 Inference in BEAST with Application to Epidemic Reconstruction. *Mol. Biol. Evol.* **37**, 1832–1842
699 (2020).
- 700 36. Bielejec, F., Lemey, P., Baele, G., Rambaut, A. & Suchard, M. A. Inferring heterogeneous
701 evolutionary processes through time: from sequence substitution to phylogeography. *Syst. Biol.* **63**,
702 493–504 (2014).
- 703 37. Cybis, G. B., Sinsheimer, J. S., Lemey, P. & Suchard, M. A. Graph hierarchies for phylogeography.
704 *Philos. Trans. R. Soc. Lond. B Biol. Sci.* **368**, 20120206 (2013).
- 705 38. Minin, V. N. & Suchard, M. A. Fast, accurate and simulation-free stochastic mapping. *Philos. Trans.*
706 *R. Soc. Lond. B Biol. Sci.* **363**, 2985–2995 (2008).
- 707 39. Rambaut, A., Drummond, A. J., Xie, D., Baele, G. & Suchard, M. A. Posterior Summarization in
708 Bayesian Phylogenetics Using Tracer 1.7. *Syst. Biol.* **67**, 901–904 (2018).

709

710 **Data availability**

711 BEAST XML input files are available at

712 https://github.com/phylogeography/SARS-CoV-2_EUR_PHYLOGEOGRAPHY

713

714 The SARS-CoV-2 genome data required for running these xmls can be downloaded from
715 <https://www.gisaid.org>. The Google COVID-19 Aggregated Mobility Research Dataset used for this study
716 is available with permission from Google LLC. The Facebook mobility data can be requested from
717 Facebook (<https://dataforgood.fb.com/>). COVID-19 incidence data was obtained from
718 <https://www.ecdc.europa.eu/en/covid-19/data>.

719

720 **Code availability**

721 The code for running BEAST analyses is available in the hmc-develop branch of the BEAST codebase
722 available at <https://github.com/beast-dev/beast-mcmc>. The tools *TreeMarkovJumpHistoryAnalyzer*,
723 *TreeStateTimeSummarizer* and *PersistenceSummarizer* are available from the master branch in the same
724 codebase.

725

726 **Acknowledgments**

727 We would like to thank all the authors who have kindly shared genome data on GISAID, and we have
728 included a table (Extended Table 6) acknowledging the authors and institutes involved.

729

730 The research leading to these results has received funding from the European Research Council under
731 the European Union's Horizon 2020 research and innovation programme (grant agreement no. 725422-
732 ReservoirDOCS) and the Bill & Melinda Gates Foundation (OPP1094793 and INV-024911). This study was
733 partially funded by EU grant 874850 MOOD and is catalogued as MOOD 005. The contents of this
734 publication are the sole responsibility of the authors and do not necessarily reflect the views of the

735 European Commission. The Artic Network receives funding from the Wellcome Trust through project
736 206298/Z/17/Z. PL acknowledges support by the Research Foundation - Flanders ('Fonds voor
737 Wetenschappelijk Onderzoek - Vlaanderen', G066215N, G0D5117N and G0B9317N). GB acknowledges
738 support from the 'Interne Fondsen KU Leuven' / Internal Funds KU Leuven under grant agreement
739 C14/18/094, and the Research Foundation – Flanders ('Fonds voor Wetenschappelijk Onderzoek -
740 Vlaanderen', G0E1420N, G098321N). MAS acknowledges support from National Institutes of Health
741 grant U19 AI135995 and R01 AI153044. SD is supported by the *Fonds National de la Recherche*
742 *Scientifique* (FNRS, Belgium). We also gratefully acknowledge support from NVIDIA Corporation with the
743 donation of parallel computing resources used for this research.

744

745 **Author contributions**

746 P.L. & S.D. designed the study, performed analyses and drafted the manuscript. V.C., C.P. and A.S.
747 provided and analyzed data. S.H., F.V., N.R., S.L. & A.T. compiled and analyzed data. A.L., B.B.O.M. and
748 M.K. contributed data. G.B. performed data analyses. M.S.G., X.J. and M.A.S. developed statistical
749 inference methodology. All authors contributed to interpreting and reviewing the manuscript.

750

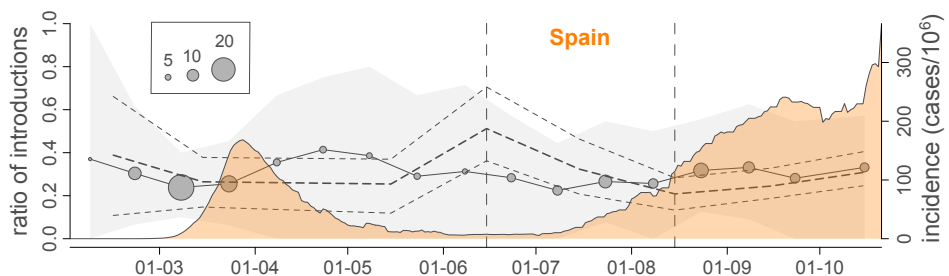
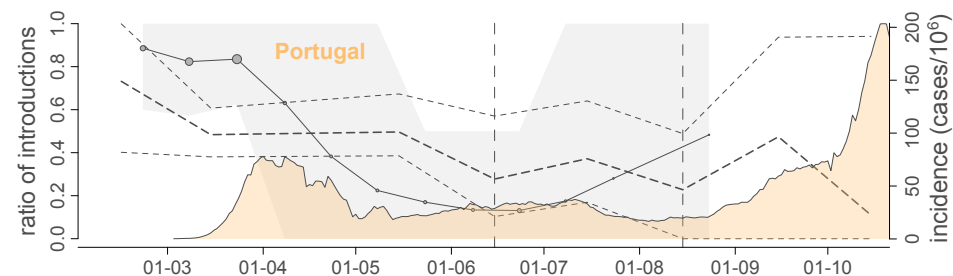
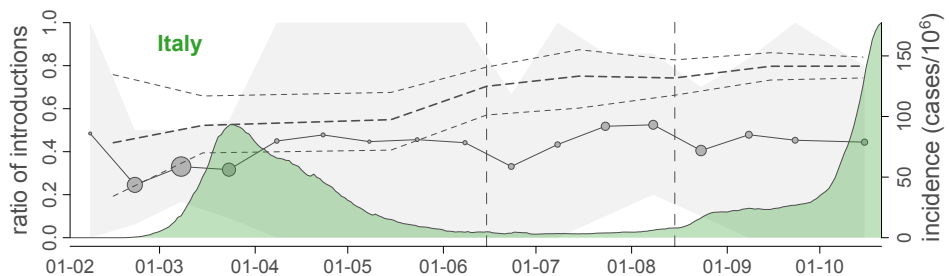
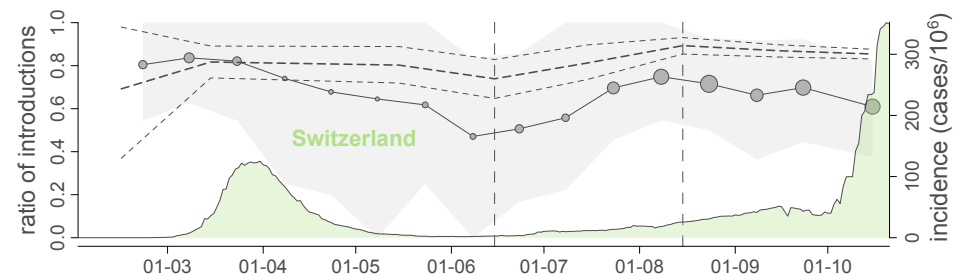
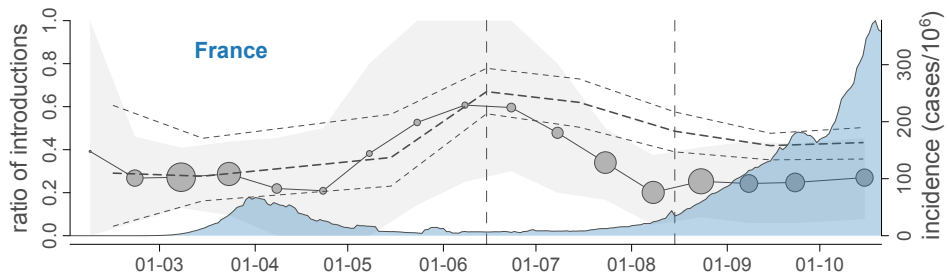
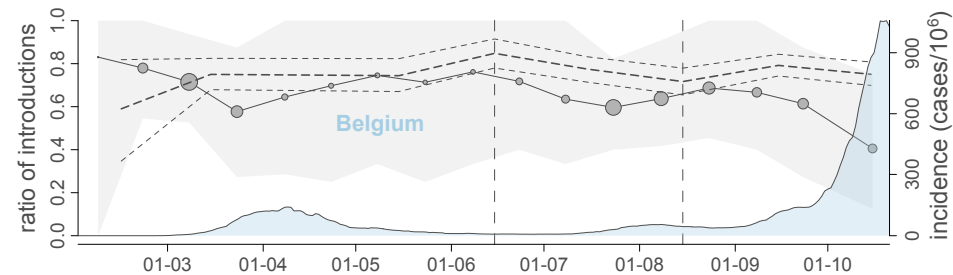
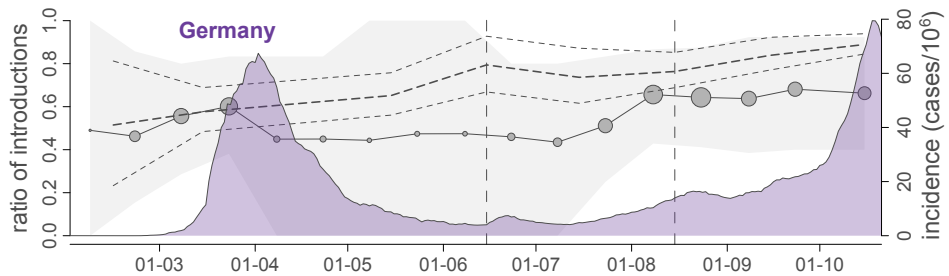
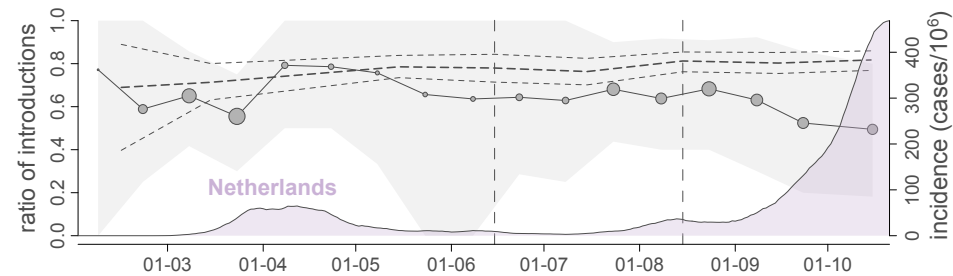
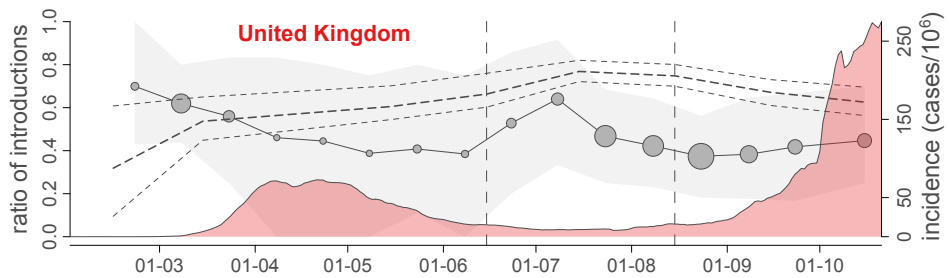
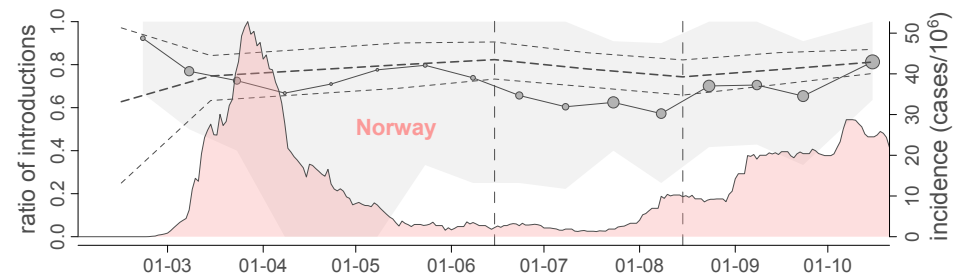
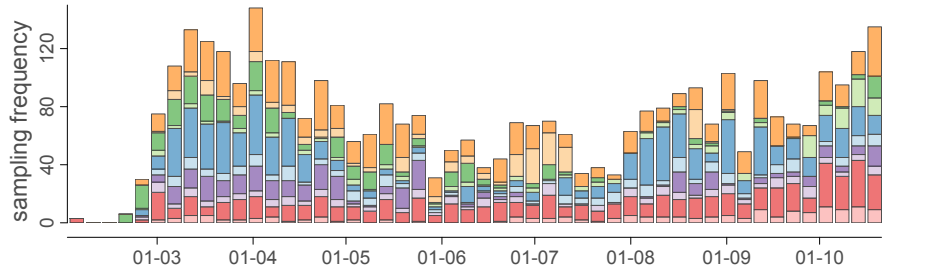
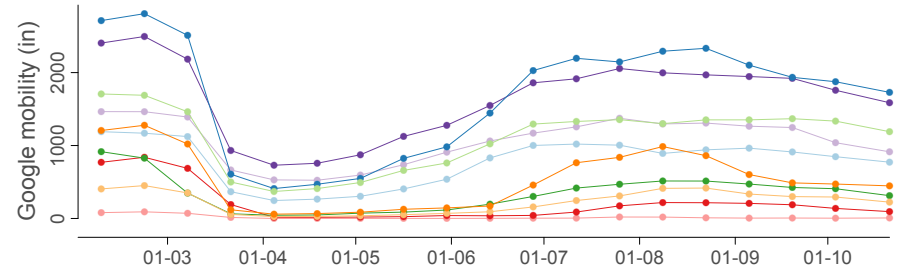
751 **Competing Interests**

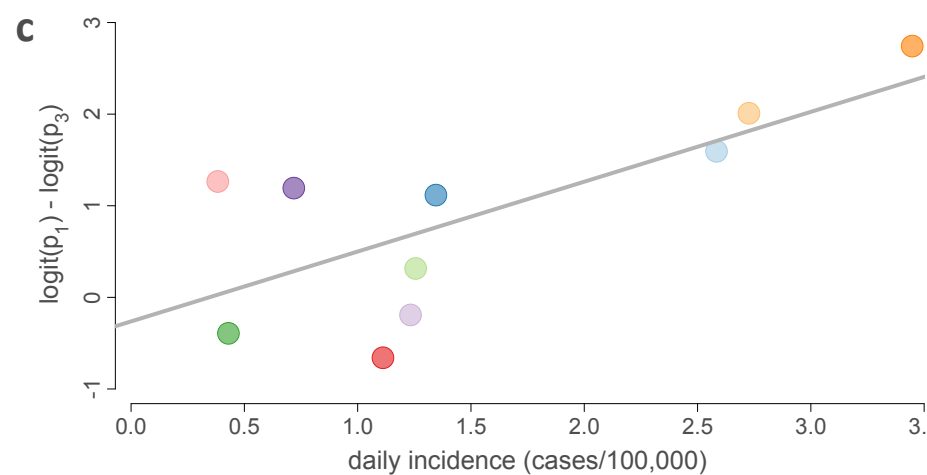
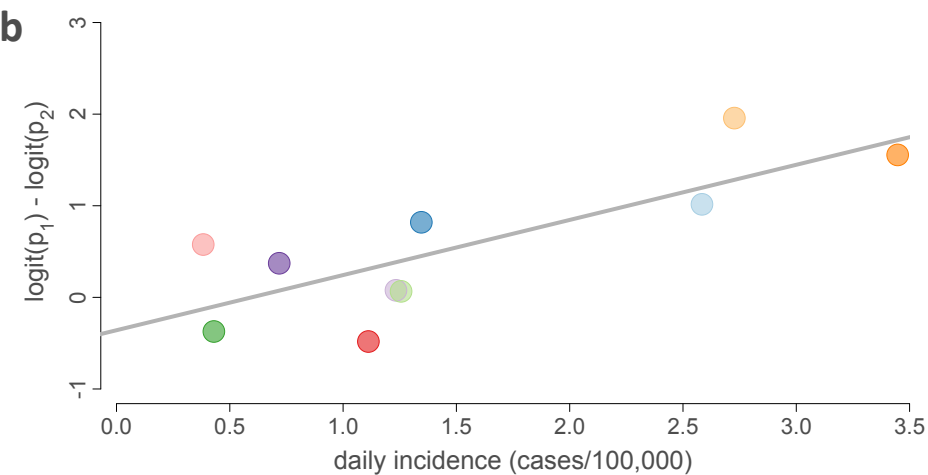
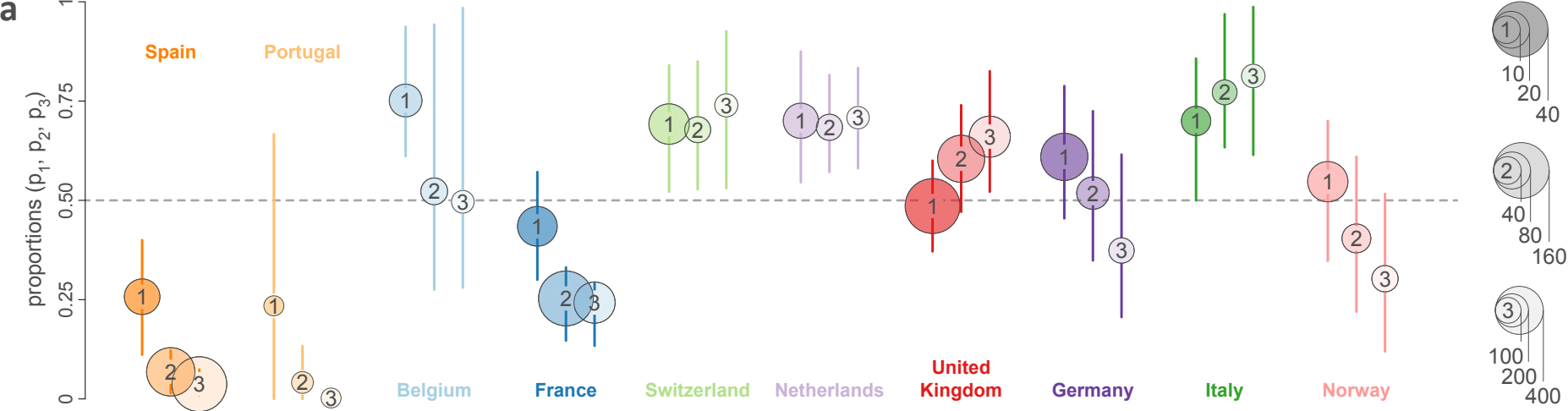
752 The authors declare no competing interests.

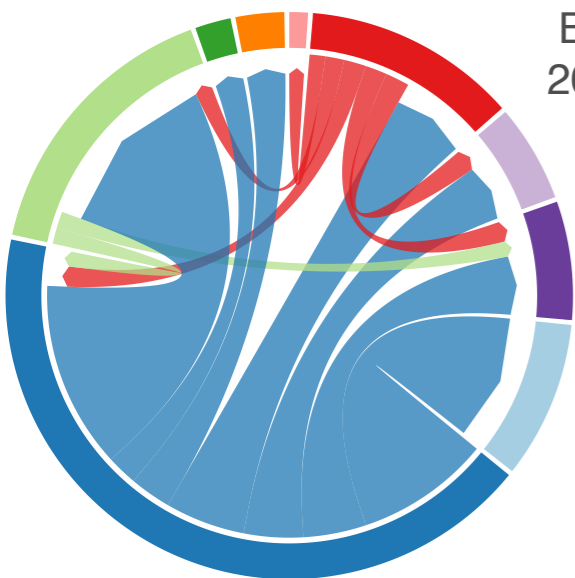
753

754 **Materials and correspondence**

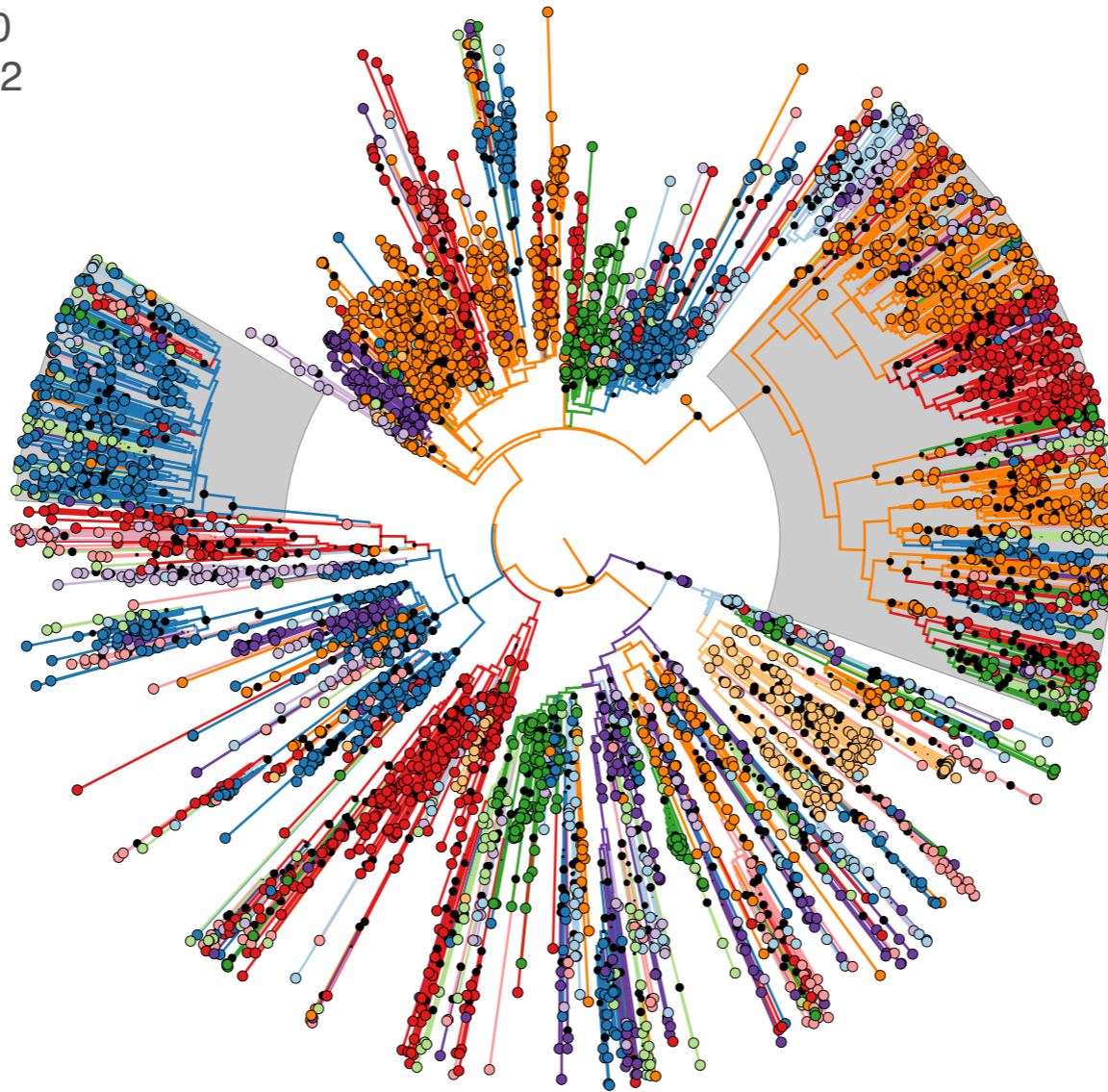
755 philippe.lemey@kuleuven.be & simon.dellicour@ulb.ac.be







B.1.160
20A.EU2



B.1.177
20A.EU1

