

1 ~~COVFlow: virus phylodynamics analyses from selected~~  
2 ~~SARS-CoV-2 sequences~~ COVFlow: performing virus  
3 phylodynamics analyses from selected SARS-CoV-2 genome  
4 sequences

5 Gonché Danesh<sup>1,\*</sup>, Corentin Boennec<sup>1</sup>, Laura Verdurme<sup>2</sup>, Mathilde Roussel<sup>2</sup>,  
Sabine Trombert-Paolantoni<sup>2</sup>, Benoit Visseaux<sup>2</sup>, Stéphanie Haim-Boukobza<sup>2</sup>, Samuel Alizon<sup>1,3</sup>

6 <sup>1</sup> MIVEGEC, CNRS, IRD, Université de Montpellier

7 <sup>2</sup> Laboratoire CERBA, Saint Ouen L'Aumône, France

8 <sup>3</sup> Center for Interdisciplinary Research in Biology (CIRB), College de France, CNRS, INSERM, Université PSL,  
9 Paris, France

10 \* corresponding author: gonche.danesh@ird.fr

## Abstract

Phylogenetic analyses can generate important and timely data to optimise public health response to SARS-CoV-2 outbreaks and epidemics. However, their implementation is hampered by the massive amount of sequence data and the difficulty to parameterise dedicated software packages. We introduce the COVFlow pipeline, accessible at <https://gitlab.in2p3.fr/ete/CoV-flow>, which allows a user to select sequences from the Global Initiative on Sharing Avian Influenza Data (GISAID) database according to user-specified criteria, to perform basic phylogenetic analyses, and to produce an XML file to be run in the **Beast2** software package. We illustrate the potential of this tool by studying two sets of sequences from the Delta variant in two French regions. This pipeline can facilitate the use of virus sequence data at the local level, for instance, to track the dynamics of a particular lineage or variant in a region of interest.

**Keywords:** COVID-19, molecular epidemiology, sequence database, phylogenetics, public health

# 1 Introduction

Millions of SARS-CoV-2 full genome sequences ~~were made available~~ have been generated since 2020 ~~through the database created by~~, and, for the majority, been made available through the Global Initiative on Sharing Avian Influenza Data (GISAID) consortium [1, 2]. This has allowed the timely monitoring of variants of concerns (VoC) with platforms such as CoVariants (CoVariants), outbreak.info [3], or CoV-Spectrum [4], and the realisation of phylogenetic analyses, e.g. via NextStrain [5].

Phylogenies represent a powerful means to analyse epidemics ~~because there is~~ with an intuitive parallel between a transmission chain and a time-scaled phylogeny of infections, which is the essence of the field known as ‘phylodynamics’ [6]. As illustrated in the case of the COVID-19 pandemic, state-of-the-art analyses allow one to investigate the spatio-temporal spread of an epidemic ~~[7]~~, superspreading events [8], and even detect differences in transmission rates between variants [9].

Phylogenetic analyses involve several technical steps to go from ~~dates~~ time-stamped virus sequence data to epidemiological parameter estimates, which can make them difficult to access to a large audience. Furthermore, the amount of data shared greatly overcomes the capacities of most software packages and imposes additional selection steps that further decrease the accessibility of these approaches. To address these limitations, we introduce the COVflow pipeline which covers all the steps from filtering the sequence data according to criteria of interest (e.g. sampling data, sampling location, virus lineage, or sequence quality) to generating a time-scaled phylogeny and an XML configuration file for a BDSKY model [10] to be run in the Beast2 software package [11].

Some pipelines already exist to assess sequence quality, filter data, infer an alignment, and infer a time-scaled phylogeny such as ~~NextClade~~ Nextclade [12] and Augur [13]. However, these do not include a ~~data filtration step based on metadata characteristics. Furthermore, performing step to~~ perform a phylogenetic analysis from the output files ~~they generate~~, which requires dedicated skills. The COVFlow pipeline addresses ~~these two limitations~~ this limitation and integrates all the steps present in separate software packages to go from the raw sequence data and metadata to the XML to be run in Beast2.

~~In this manuscript~~ Here, we present the architecture of the pipeline and apply it to data from the French epidemic, which has been poorly analysed (but see [14–16]). Focusing on sequences belonging

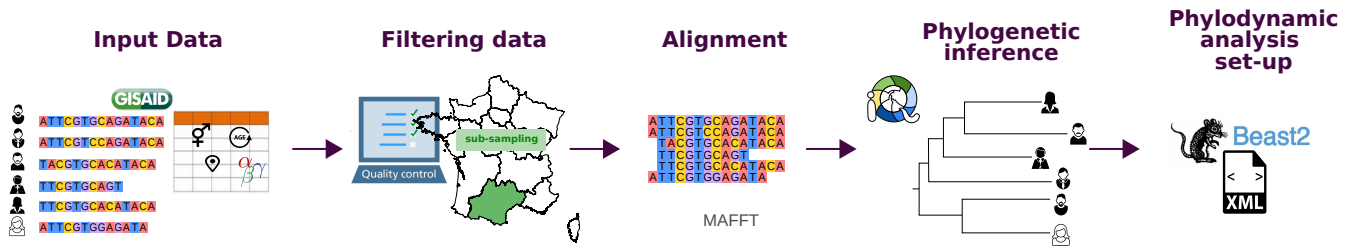


Figure 1: **Structure of the COVFlow pipeline.** The input data correspond to FASTA sequences and metadata provided by the GISAID. The data filtering is done using a YAML configuration file. The sequence alignment is performed with [MAFFT](#) and the phylogenetic inference with [IQ-TREE](#). The pipeline generates an XML file that can be directly used with [Beast2](#).

51 to the Delta variant collected in France in two regions, Ile-de-France, and Provence-Alpes-Cote-d’Azur,  
 52 by a specific French laboratory (CERBA), we illustrate the [pipeline](#) accessibility, flexibility, and public  
 53 health relevance of the [COVFlow pipeline](#).

## 54 2 Methods

55 COVFlow is a bioinformatics pipeline for phylogenetic and phylodynamic analysis of SARS-CoV-2  
 56 genome sequences. It is based on the Snakemake workflow management system [\[17\]](#) and its de-  
 57 pendencies are easily installed via a conda virtual environment. Snakemake ensures reproducibility,  
 58 while Conda (<https://docs.conda.io/en/latest/>) and Bioconda [18] [allow](#) [allows](#) for version con-  
 59 trol of the programs used in the pipeline. Overall, the pipeline is easy to install and avoids dependency  
 60 conflicts.

### 61 Pipeline configuration

62 The pipeline workflow is configured using a [configuration file, in a YAML format.](#) ~~The configuration~~  
 63 ~~file~~ [YAML configuration file, which](#) must contain the path to the sequence data file, the path to the  
 64 metadata file, and the prefix chosen for the output files. Each parameter of the pipeline following steps  
 65 has a default value, which can be modified by the user in the [YAML](#) configuration file.

## 2.1 ~~Input data~~

### Input data

The input data analysed by COVFlow are sequence data and metadata, corresponding to patient properties, that can be downloaded from the Global Initiative on Sharing All Influenza Data (GISAID, <https://www.gisaid.org/>). The sequence data are in a FASTA format file. The metadata downloaded contains details regarding the patient’s sequence ID (column named ‘strain’), the sampling dates (column ‘date’), the region, country, and division where the sampling ~~has been~~ was made (columns respectively named ‘region’, ‘country’, and ‘division’). It also lists the virus lineage assigned by the Pangolin tool [19], and the age and sex of the patient (columns respectively named ‘pango\_lineage’, ‘age’, and ‘sex’).

### Data filtering

The first step implemented in the pipeline performs quality filtering. By default, genomic sequences that are shorter than 27,000 bp, or that have more than 3,000 missing data (i.e. N bases) and more than 15 non-ATGCN bases are excluded. These parameter values can be modified by the user. Sequences belonging to non-human or unknown hosts are also excluded. Sequences for which the sampling date is more recent than the submission date, or for which the sampling date is unclear (e.g. missing day) are also excluded. Finally, duplicated sequences and sequences that are flagged by the Nextclade tool [12] with an overall bad quality (Nextclade QC overall status ‘bad’ or ‘mediocre’) are also removed.

The sequence data is then further filtered following the user’s criteria. These include Pangolin lineages, sampling locations (regions, countries, or divisions), and sampling dates. In addition to specifying the maximum and/or minimum sampling dates, the user can specify a sub-sampling scheme of the data with a number or percentage of the data per location and/or per month. For example, the user can decide to keep  $x\%$  of the data per country per division per month or to keep  $y$  sequence data per division. Finally, more specific constraints can be given using a JSON format file with three possible actions: i) keep only rows (i.e. sequences) that match or contain a certain value, ii) remove rows that match or contain a certain value, and iii) replace the value of a column by another value for specific rows with a column that matches or contains a certain value. The last action can be used

93 to correct the metadata, for instance, if the division field is not filled in but can be inferred from the  
94 names of the submitting laboratory. The JSON file can be composed of multiple key-value pairs, each  
95 belonging to one of the three actions. For example, the user can specify to keep only male patients and  
96 to remove data from one particular division while setting the division of all the samples submitted by  
97 a public hospital from the Paris area (i.e. the APHP) to the value 'Ile-de-France'.

## 98 **Aligning and masking**

99 The set of sequences resulting from the data filtering is then divided into temporary FASTA files with  
100 a maximum number of 200 sequences per file. For each subset, sequences are aligned to the reference  
101 genome MN908947.3 using ~~MAFFT~~ MAFFT v7.305 [20] with the 'keeplength' and '~~addfragments~~add'  
102 options. All the aligned sequences are then aggregated into a single file. Following earlier studies,  
103 the first 55~~and~~, the last 100 sites, and other sites recommended from [https://github.com/W-L/](https://github.com/W-L/ProblematicSites_SARS-CoV2)  
104 [ProblematicSites\\_SARS-CoV2](https://github.com/W-L/ProblematicSites_SARS-CoV2) of the alignment are then masked to improve phylogenetic inference  
105 (<http://virological.org/t/issues-with-sars-cov-2-sequencing-data/473>).

## 106 **Inferring and time-scaling a phylogeny**

107 A maximum-likelihood phylogenetic tree is estimated using ~~IQ-TREE~~ IQ-TREE v2.1.2 [21] under a  
108 GTR substitution model from the alignment. The resulting phylogeny is time-scaled using ~~TreeTime~~  
109 TreeTime v0.8.1 [22]. By default, the tree is rooted using two ancestral sequences (Genbank accession  
110 numbers MN908947.3 and MT019529.1) as an outgroup, which is then removed, with a fixed clock rate  
111 of  $8 \cdot 10^{-4}$  substitutions per position per year [23] and a standard deviation of the given clock rate of  
112 0.0004. These parameters can be modified by the user. The output phylogeny is in a Newick format  
113 file.

## 114 **BDSKY XML file generating for BEAST 2**

115 The Bayesian birth-death skyline plot (or BDSKY) method allows the inference of the effective ~~reproductive~~  
116 reproduction number from genetic data or directly from a phylogenetic tree, by estimating transmis-  
117 sion, recovery, and sampling rates [10]. This method allows these parameters to vary through time and

118 is implemented within the ~~BEAST2 software framework~~ BEAST 2 software package [11].

119 Performing a BDSKY analysis requires setting an XML file specifying the parameters for the priors.  
120 As in any Bayesian analysis, this step is extremely important. The default settings in BEAST2 have  
121 been chosen to minimise the risk of errors. COVFlow builds on most of these with some modifications  
122 to fit the needs of large SARS-CoV-2 phylogenies.

123 The most important change has to do with the inference of the phylogeny. This can be done by  
124 ~~BEAST2~~ BEAST 2 but to minimise computation speed and allow for the analysis of large phylogenies,  
125 the pipeline sets the time-scaled phylogeny from the previous step in the XML file.

126 The default XML file assumes that there are two varying effective ~~reproductive~~ reproduction  
127 numbers to estimate, with a lognormal prior distribution,  $\text{LogNorm}(M = 0, S = 1)$ , resulting in a median  
128 of 1, the 95% quantiles falling between 7.10 and 0.14, and a starting value of ~~±1.0~~. This prior is  
129 adapted to such virus epidemics and, as we will see below, can be edited if needed. The default prior  
130 for the rate of end of the infectious period is a uniform distribution,  $\text{Uniform}(10, 300)$ , resulting in a  
131 median of 155[17.3; 293]years<sup>-1</sup>, with a starting value of 70 years<sup>-1</sup>, and is assumed to be constant  
132 over time. ~~This~~ The inverse of the rate of end of the infectious period is the average infectious period.  
133 This default prior yields infectious periods varying from ~~1.2 to 36.5 days~~ 0.034 year (1.2 days) to 0.1  
134 year (36.5 days), which is ~~relevant~~ consistent with the biology of SARS-CoV-2 infections [24]. Usually,  
135 little or no sampling effort is made before the first sample was collected. Therefore, by default, we  
136 assume two sampling proportions: before the first sampling date it is set to zero, and after the default  
137 prior is a beta distribution,  $\text{Beta}(\alpha = 1, \beta = 1)$ , with a starting value of 0.01, translating in a median  
138 of 0.50 ([0.025; 0.975]). The non-zero sampling proportion is assumed to remain constant during the  
139 time the samples were collected. The method can also estimate the date of origin of the index case,  
140 which corresponds to the total duration of the epidemic. Since the tree is assumed to be a sampled  
141 tree, and not a complete one, the origin is always earlier than the time to the most recent common  
142 ancestor of the tree. Hence, the prior distribution's starting value and upper value must be higher  
143 than the tree height. This condition is always checked when running the pipeline. The default prior  
144 for this parameter prior is a uniform distribution ~~Uniform(0, 2) years~~ Uniform(0, height + 2) years, with  
145 a starting value of height, with height as the maximum height of the inferred time-scaled tree.

Note that, although the default priors are designed to minimise the risk of bias in the results and the pipeline checks for the origin parameter prior, the choice of the priors is essential and may impact the phylodynamic inference of parameters.

In the COVFlow configuration file, the user can modify the distribution shapes, the starting values, the upper and lower values, and the dimensions for each of these parameters to estimate, and set the dates at which the parameter estimation changes. The length of the MCMC chain and the sampling frequency, which are by default set to 10,000,000 and 100,000 respectively, can also be modified.

The BEAST2-BEAST 2 inference itself is not included in the pipeline. The reason for this is that a preliminary step (i.e. installing the BDSKY package) needs to be performed by the user. Similarly, the analysis of the BEAST2 output log files needs to be performed by the user via Tracer [25] or a dedicated R script available on the COVflow Gitlab page.

### 3 Results

~~Analysing the SARS-CoV-2 Delta variant epidemics in French regions using the COVFlow pipeline.~~ a) Geographical sub-sampling using at most 50 sequences per month for the Delta variant in Ile-de-France (IdF, in red), Provence-Alpes-Côte d'Azur (PACA, in green), and in all of France collected by CERBA laboratory. b) Time-scaled phylogenies generated using sub-sampled data from IdF (in red), PACA (green), and all of France (in orange). c) Temporal variations of the effective reproductive number ( $R_e$ ) of the Delta variant in IdF (red), in PACA (green), and France (orange) estimated using Beast2. The last panel was generated using Beast2. In panel c, the lined show the median values and the shaded area the 95% highest posterior density.

We illustrate the potential of the ~~Compared with the COVFlow~~Nextstrain pipeline by performing a phylodynamic analysis of a specific COVID-19 lineage, here the Delta variant (Pango lineage B.1.617.2), in two regions of a country, here Ile-de-France and Provence-Alpes-Côte d'Azur in France (Figure 2(a)) [5]. COVflow allows a more flexible filtering stage using the JSON file. Furthermore, the sub-sampling can either be based on the number of data points or on the percentage of available data and the latter option is currently not possible with Nextstrain. The masking sites strategy is also different between the two pipelines. Finally, and perhaps most importantly, COVflow configures an XML file for a BDSKY



173 phylogenetic analysis in Beast 2, allowing for more detailed phylogenetic analyses.

174 ~~We downloaded~~

## 175 Illustration study with French data

176 We applied COVFlow to analyse GISAID data by downloading sequence data and metadata from the  
177 GISAID platform for the GK clade corresponding to the lineage B.1.617.2 available on ~~the~~ April 22,  
178 2022, which amounted to 4,212,049 sequences. Using the pipeline and the editing of ~~the~~ its JSON  
179 file, we cleaned the sequence data, selected the data collected by a specific large French laboratory  
180 (CERBA), selected the data from two regions of interest (the Ile-de-France region for ~~the first analysis~~  
181 first analysis, and the Provence-Alpes-Côte d'Azur region for ~~the a~~ second analysis), and sub-sampled  
182 the data to keep up to 50 sequences per month. ~~For the third analysis, which~~ These two regions were  
183 chosen because they had some of the highest coverage in the dataset, while being in different parts of  
184 France. Our third analysis included the whole country, ~~so~~ we sub-sampled the data to keep up to 50  
185 sequences per month per French region. ~~This~~ The other parameters of the pipeline were default except  
186 for the number of windows for the effective reproduction numbers in the BDSKY analysis which was  
187 set to 9 with a change-point time every month from June 01, 2021, to January 01, 2022.

188 To evaluate the robustness of the inference, we performed 5 independent COVFlow runs for France,  
189 using the pipeline configuration described above for the France analysis. For each run, we manually  
190 extracted the two major clades representing at least 20% of the leaves from the resulting phylogeny and  
191 used a Python script of the COVFlow pipeline to generate two XML files. For each BDSKY analysis,  
192 9 effective reproduction numbers were estimated over the same time periods.

193 To assess the validity of the BDSKY results, we extracted SARS-CoV-2 PCR screening data from  
194 <https://www.data.gouv.fr/fr/datasets/r/5c4e1452-3850-4b59-b11c-3dd51d7fb8b5>. More precisely,  
195 we used the positivity rate at the national level and in the two regions of interest. The effective  
196 reproduction number ( $R_e$ ) was estimated using the EpiEstim R package [26, 27]. The data were  
197 smoothed out using a 7-days rolling average, in order to compensate for the reporting delays.

198 The files necessary to generate these results are provided in Appendix.

### 3 Results

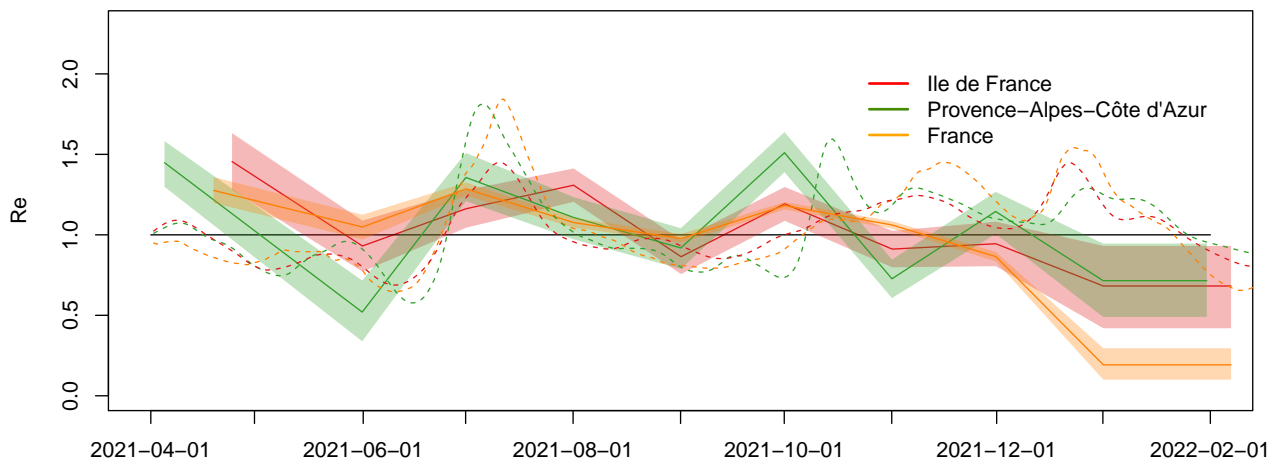
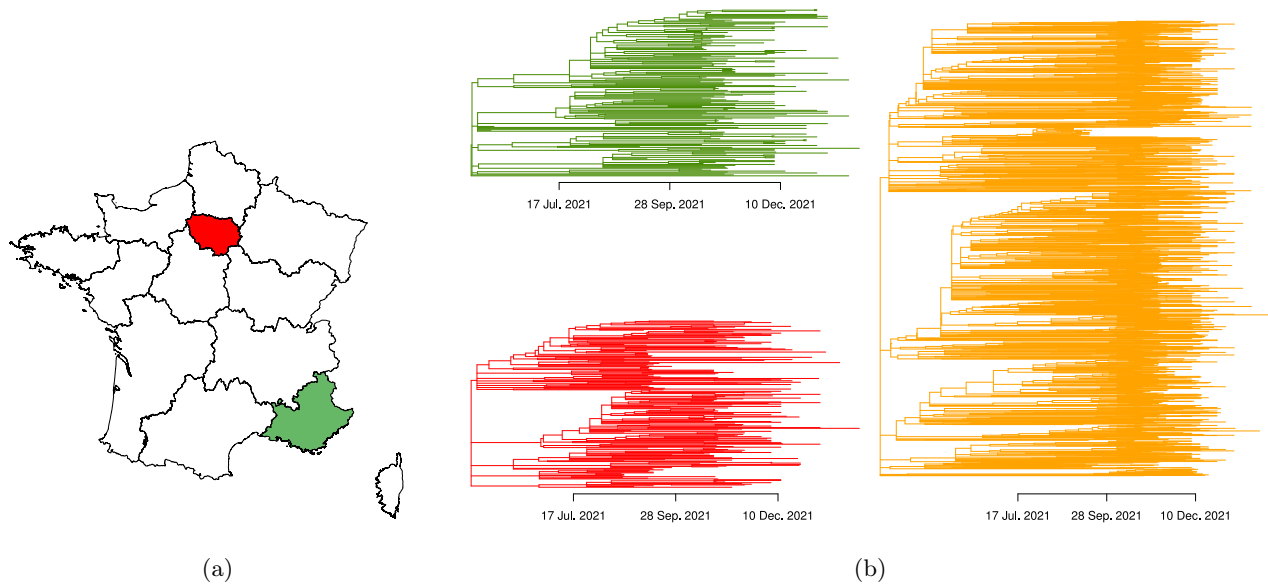
We illustrate the potential of the COVFlow pipeline by performing a phylodynamic analysis of a specific COVID-19 lineage, here the Delta variant (Pango lineage B.1.617.2), in two regions of a country, here Ile-de-France and Provence-Alpes-Côte d’Azur in France (Figure 2(a)).

The COVflow runs resulted in the selection of 176 SARS-CoV-2 genomes for Ile-de-France (IdF), 221 genomes for Provence-Alpes-Côte d’Azur (PACA), and 1,575 genomes for France. In this example, the two regions were chosen because they had some of the highest coverage in the dataset, while being in different parts of France. The other parameters of the pipeline were default except that for the number of windows for the effective reproductive numbers in the BDSKY analysis which was set to 10.

The first output of the pipeline is the time-scaled phylogeny inferred from the sequences. In Figure 2(b), we show the one for each of the two French regions considered and the one for France the whole country. This already allow allows us to visualise the date of origin of the epidemic associated with the sequences sampled. Furthermore More generally, the shape of the phylogeny phylogenies can reflect the number of introductions epidemic spread in the locality studied, e.g. the number of external introductions.

The second output of the pipeline is the XML file for a BDSKY model that can be run into Beast2. In Figure 2(c), we show the temporal variations in the effective reproduction number ( $R_e$ ), that which is the average number of secondary infections caused by an infected individual at a given date. If  $R_e < 1$ , the epidemic is decreasing and if  $R_e > 1$  it is growing.

The results show that the importation of the Delta variant epidemic seems to occurred earlier and more frequently have started earlier in PACA than in IdF in early 2021. Furthermore, in early July 2021, we see that the epidemic wave started in PACA before IdF. This In both regions (and in France), the growth of the Delta variant in June is consistent with previous results showing the transmission advantage of 79% over the Alpha variant during this time period [28]. Furthermore, the earlier start in PACA is consistent with the beginning of the school holidays and, PACA being a densely populated region in the summer. During the summer 2021, the epidemic growth in these two regions, IdF and Note that IdF, as PACA, was more important than above the French average. In the fall, which is



(c)

Figure 2: Analysing the SARS-CoV-2 Delta variant epidemics in French regions using the COVFlow pipeline. a) Geographical sub-sampling using at most 50 sequences per month for the Delta variant in Ile-de-France (IdF, in red), Provence-Alpes-Côte d'Azur (PACA, in green), and in all of France collected by CERBA laboratory. b) Time-scaled phylogenies generated using sub-sampled data from IdF (in red), PACA (green), and all of France (in orange). c) Temporal variations of the effective reproduction number ( $R_e$ ) of the Delta variant in IdF (red), in PACA (green), and France (orange) estimated using Beast2 from phylogenies in solid lines, and estimated using Epiestim from incidence data in dashed lines. The last panel was generated using Beast2. In panel c, the solid lines show the median values and the shaded area the 95% highest posterior density.

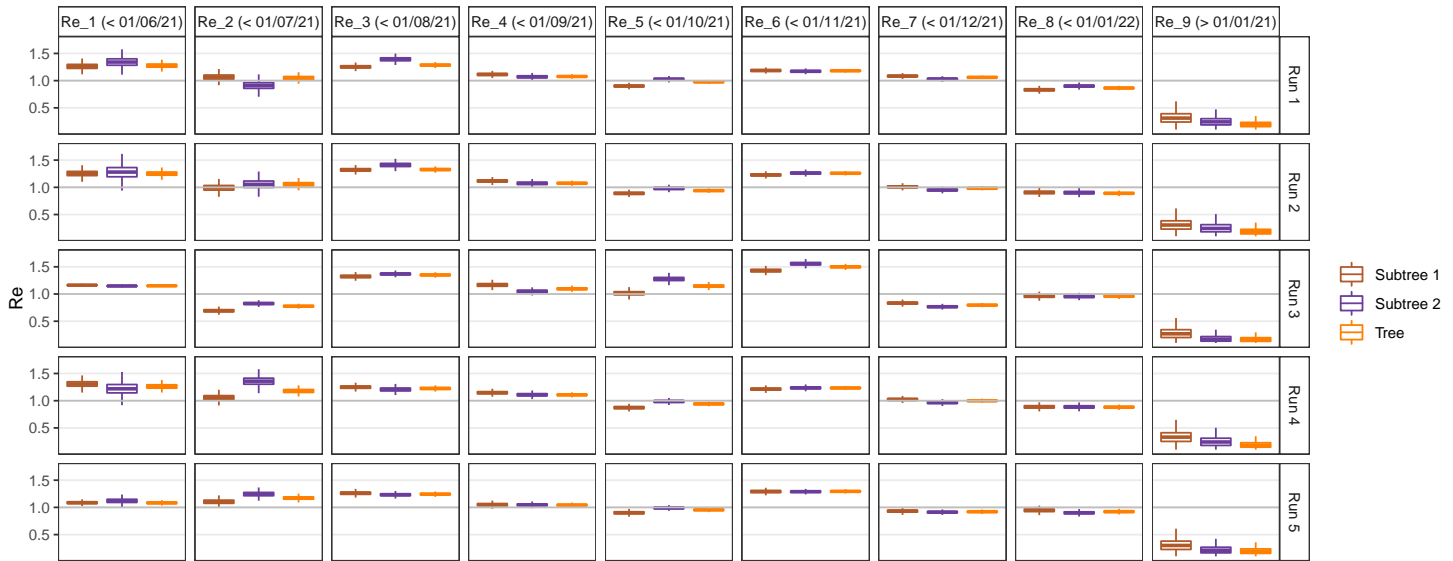


Figure 3: Estimations of the effective reproduction number  $R_e$  of the Delta variant in France for 5 different COVflow runs. For each run, the  $R_e$  were estimated from the inferred phylogenetic tree, and for the two principal clades, denoted Subtree 1 and Subtree 2. For each tree, 9 different  $R_e$  were estimated, with a changing point date every month from 2021-06-01 to 2022-01-01.

227 also unsurprising given the density and international connections of the region.

228 Early in the fall of 2021, the back-to-school period led to an epidemic rebound in France. The  
 229 associated epidemic growth was again stronger and earlier in PACA than in IdF. Furthermore, given  
 230 the superposition of the French and PACA  $R_e$  curves, it is possible that PACA drove the national  
 231 epidemic at the time. Finally, we see that the national average. Furthermore, contrarily to IdF or  
 232 France, PACA experience experienced a period of Delta variant growth following the winter holidays.  
 233 These are more difficult to explain but could be linked to local differences in terms of behaviour.

234 Finally, we see a clear slowdown in the Delta variant epidemic at the end of 2021. This final result is  
 235 consistent with large French surveillance data obtained through variant-specific PCR tests that show a  
 236 high proportion of infections with the S:L452R mutation in PACA compared to the other French regions  
 237 during the Omicron wave is likely linked to the extension of the 3<sup>rd</sup> vaccination dose, to changes in  
 238 French behavior, but also to emergence of Omicron BA.1 variant, which was shown to have a growth  
 239 advantage over the Delta variant [29].

240 When comparing BDSKY estimates with that of EpiEstim on the screening tests (dashed lines), we  
 241 generally found consistent results. However, we did observe a shift in  $R_e$  peaks. This is consistent with

242 the fact that methods based on incidence have an intrinsic delay due to the lag between the date of  
243 the infection and that of the PCR testing. For phylodynamics, this delay is, in theory, less important  
244 since the methods focus on virus evolution. EpiEstim estimates detect an epidemic growth in IdF and  
245 then PACA at the end of 2021 but this is expected because PCR tests do not discriminate between  
246 lineages and the end of the 2021 year saw the rise of the Omicron BA.1 variant [29].

247 Finally, a worry with phylodynamics is that the results depend on the sequences chosen. Moreover,  
248 considering the whole phylogeny incorporates importation events that are not included explicitly in  
249 the underlying birth-death model assumed by the BDSKY methods. Fig. 3, we show that the effective  
250 reproduction numbers estimated from the main subtrees are quantitatively similar to the  $R_e$  estimated  
251 from the whole phylogenetic tree. Furthermore, the estimations are all similar for different runs  
252 suggesting that the BDSKY framework is robust to phylogenetic tree uncertainty.

## 253 4 Discussion

254 The COVID-19 pandemic constitutes a qualitative shift in terms of the generation, sharing, and analysis  
255 of virus genomic sequence data. The GISAID initiative allowed the rapid sharing of SARS-CoV-2  
256 sequence data, which is instrumental for local, national, and international public health structures that  
257 need to provide timely reports on the sanitary situation. At a more fundamental level, this genomic  
258 data is also key to furthering our understanding of the spread and evolution of the COVID-19 pandemic  
259 [30], especially in low-resource countries [31].

260 We ~~elaborate~~elaborated the COV-flow pipeline, which allows users to perform all the steps from  
261 raw sequence data to phylodynamics analyses. In particular, it can select sequences from the GISAID  
262 ~~datased~~dataset based on metadata, perform a quality check, align the sequences, infer a phylogeny,  
263 root this phylogeny into time, and generate an XML file for Beast2 analysis (we also provide scripts to  
264 analyse the outputs). Furthermore, COV-flow can also readily allow the implementation of subsampling  
265 schemes per location and per date. This can help balance the dataset and also be extremely useful to  
266 perform sensitivity analyses and explore the robustness of the phylodynamic results.

267 A future extension ~~will~~could consist in including other Beast2 population dynamics models, for  
268 instance, the Bayesian Skyline model, which is not informative about  $R_0$  but is potentially less sensitive

269 to variations in sampling intensity ~~as it assumes sampling is negligible~~. Another extension ~~will~~ could be  
270 to use other databases to import SARS-CoV-2 genome data, e.g. that published by NCBI, via LAPIS  
271 (Lightweight API for Sequences).

272 Beast2 can simultaneously infer population dynamics parameters and phylogenies, which is an  
273 accurate way to factor in phylogenetic uncertainty [11]. However, this global inference is particu-  
274 larly computationally heavy and is out of reach for large data sets. To circumvent this problem,  
275 we perform the phylogenetic inference first using less accurate software packages and then impose  
276 the resulting phylogeny into the Beast2 XML file. An extension of the pipeline could offer the  
277 user to also perform the phylogenetic inference, for instance by using the so-called ‘Thorney Beast’  
278 ([https://beast.community/thorney\\_beast](https://beast.community/thorney_beast)) implemented in Beast 1.10 [32].

279 Finally, it is important to stress that phylogenetic analyses are always dependent on the sampling  
280 scheme [33–36]. If most of the sequences come from contact tracing in dense clusters, the analysis will  
281 tend to overestimate epidemic spread. This potential bias can be amplified by the sequence selection  
282 feature introduced in the pipeline. An advantage of COVFlow is that it can perform spatio-temporal  
283 subsampling but additional studies are needed to identify which are the most appropriate subsampling  
284 schemes to implement.

## 285 Acknowledgement

286 The authors acknowledge further support from the CNRS, the IRD and the ~~itrop~~ i-Trop HPC (South  
287 Green Platform) at IRD ~~montpellier~~ Montpellier, which provided HPC resources that contributed to  
288 the results reported here (<https://bioinfo.ird.fr/>).

289 The authors thank the Experimental and Theoretical Evolution team from Maladies Infectieuses et  
290 Vecteurs: Écologie, Génétique, Évolution et Contrôle, University of Montpellier, for discussion, as well  
291 as the EMERGEN consortium (complete member list in Supplementary Materials).

292 This project was supported by the Agence Nationale de la Recherche Maladies Infectieuses Émer-  
293 gentes to the MODVAR project (grant no. ANRS0151).

## 294 Authors contributions

295 GD and SA conceived the study, GD built the pipeline and performed the analyses, CB contributed to  
296 the implementation of the pipeline, LV, MR, STP, BV, and SHB contributed genetic sequence data,  
297 SA and GD wrote a first version of the manuscript.

## 298 Data and scripts

299 The sequences analysed were generated by CERBA and uploaded to GISAID.

300 The R scripts, along with all the files generated by the pipeline and used for the analyses (XML  
301 files, FASTA alignments, time-scaled phylogenies) are provided in Supplementary Materials.

302 The pipeline itself can be accessed on the Git public repository [https://gitlab.in2p3.fr/ete/](https://gitlab.in2p3.fr/ete/CoV-flow)  
303 CoV-flow

## 304 Conflict of Interest

305 The authors of this preprint declare that they have no financial conflict of interest with the content of  
306 this article.

## 307 References

- 308 [1] Elbe, S. & Buckland-Merrett, G., 2017 Data, disease and diplomacy: GISAID’s innovative contri-  
309 bution to global health. *Global Challenges* **1**, 33–46. (doi:<https://doi.org/10.1002/gch2.1018>).
- 310 [2] Khare, S., Gurry, C., Freitas, L., Schultz, M. B., Bach, G., Diallo, A., Akite, N., Ho, J., Lee, R. T.,  
311 Yeo, W. *et al.*, 2021 GISAID’s Role in Pandemic Response. *China CDC Weekly* **3**, 1049–1051.  
312 (doi:[10.46234/ccdcw2021.255](https://doi.org/10.46234/ccdcw2021.255)).
- 313 [3] Latif, A. A., Mullen, J. L., Alkuzweny, M., Tsueng, G., Cano, M., Haag, E., Zhou, J., Zeller, M.,  
314 Matteson, N., Wu, C. *et al.*, 2021. outbreak.info: Lineage comparison.
- 315 [4] Chen, C., Nadeau, S., Yared, M., Voinov, P., Xie, N., Roemer, C. & Stadler, T., 2022 CoV-  
316 Spectrum: analysis of globally shared SARS-CoV-2 data to identify and characterize new variants.  
317 *Bioinformatics* **38**, 1735–1737. (doi:[10.1093/bioinformatics/btab856](https://doi.org/10.1093/bioinformatics/btab856)).

- 318 [5] Hadfield, J., Megill, C., Bell, S. M., Huddleston, J., Potter, B., Callender, C., Sagulenko, P., Bed-  
319 ford, T. & Neher, R. A., 2018 Nextstrain: real-time tracking of pathogen evolution. *Bioinformatics*  
320 **34**, 4121–4123. (doi:10.1093/bioinformatics/bty407).
- 321 [6] Grenfell, B. T., Pybus, O. G., Gog, J. R., Wood, J. L., Daly, J. M., Mumford, J. A. & Holmes,  
322 E. C., 2004 Unifying the epidemiological and evolutionary dynamics of pathogens. *Science* **303**,  
323 327–32. (doi:10.1126/science.1090727).
- 324 [7] Plessis, L. d., McCrone, J. T., Zarebski, A. E., Hill, V., Ruis, C., Gutierrez, B., Raghwani, J.,  
325 Ashworth, J., Colquhoun, R., Connor, T. R. *et al.*, 2021 Establishment and lineage dynamics of  
326 the SARS-CoV-2 epidemic in the UK. *Science* (doi:10.1126/science.abf2946).
- 327 [8] Alizon, S., 2021 Superspreading genomes. *Science* **371**, 574–575. (doi:10.1126/science.abg0100).
- 328 [9] Volz, E., Hill, V., McCrone, J. T., Price, A., Jorgensen, D., O’Toole, , Southgate, J.,  
329 Johnson, R., Jackson, B., Nascimento, F. F. *et al.*, 2021 Evaluating the Effects of SARS-  
330 CoV-2 Spike Mutation D614G on Transmissibility and Pathogenicity. *Cell* **184**, 64–75.e11.  
331 (doi:10.1016/j.cell.2020.11.020).
- 332 [10] Stadler, T., Kühnert, D., Bonhoeffer, S. & Drummond, A. J., 2013 Birth-death skyline plot reveals  
333 temporal changes of epidemic spread in HIV and hepatitis C virus (HCV). *Proc Natl Acad Sci*  
334 *USA* **110**, 228–33. (doi:10.1073/pnas.1207965110).
- 335 [11] Bouckaert, R., Heled, J., Kühnert, D., Vaughan, T., Wu, C.-H. *et al.*, 2014 Beast 2: a software  
336 platform for bayesian evolutionary analysis. *PLoS Comput Biol* **10**, e1003537.
- 337 [12] Aksamentov, I., Roemer, C., Hodcroft, E. B. & Neher, R. A., 2021 Nextclade: clade assignment,  
338 mutation calling and quality control for viral genomes. *Journal of Open Source Software* **6**, 3773.  
339 (doi:10.21105/joss.03773).
- 340 [13] Huddleston, J., Hadfield, J., Sibley, T. R., Lee, J., Fay, K., Ilcisin, M., Harkins, E., Bedford, T.,  
341 Neher, R. A. & Hodcroft, E. B., 2021 Augur: a bioinformatics toolkit for phylogenetic analyses of  
342 human pathogens. *Journal of Open Source Software* **6**, 2906. (doi:10.21105/joss.02906).



- 343 [14] Danesh, G., Elie, B., Michalakis, Y., Sofonea, M. T., Bal, A., Behillil, S., Destras, G., Boutolleau,  
344 D., Burrel, S., Marcelin, A.-G. *et al.*, 2021 Early phylogenetics analysis of the COVID-19 epidemic  
345 in France. *Peer Community Journal* **1**, e45. (doi:10.24072/pcjournal.40).
- 346 [15] Gambaro, F., Baidaliuk, A., Behillil, S., Donati, F., Albert, M., Alexandru, A., Vanpeene, M.,  
347 Bizard, M., Brisebarre, A., Barbet, M. *et al.*, 2020 Introductions and early spread of SARS-CoV-2  
348 in France. *Eurosurveillance* **25**, 2001200. (doi:10.2807/1560-7917.ES.2020.25.26.2001200).
- 349 [16] Coppée, R., Blanquart, F., Jary, A., Leducq, V., Ferré, V. M., Franco Yusti, A. M., Daniel,  
350 L., Charpentier, C., Lebourgeois, S., Zafilaza, K. *et al.*, 2023 Phylogenetics of SARS-CoV-2 in  
351 France, Europe, and the world in 2020. *eLife* **12**. (doi:10.7554/eLife.82538).
- 352 [17] Köster, J. & Rahmann, S., 2012 Snakemake—a scalable bioinformatics workflow engine. *Bioin-*  
353 *formatics* **28**, 2520–2522. ISSN 1367-4803. (doi:10.1093/bioinformatics/bts480).
- 354 [18] Grüning, B., Dale, R., Sjödin, A., Chapman, B. A., Rowe, J., Tomkins-Tinch, C. H., Valieris,  
355 R. & Köster, J., 2018 Bioconda: sustainable and comprehensive software distribution for the life  
356 sciences. *Nature methods* **15**, 475–476. (doi:10.1038/s41592-018-0046-7).
- 357 [19] O’Toole, , Scher, E., Underwood, A., Jackson, B., Hill, V., McCrone, J. T., Colquhoun, R., Ruis,  
358 C., Abu-Dahab, K., Taylor, B. *et al.*, 2021 Assignment of epidemiological lineages in an emerging  
359 pandemic using the pangolin tool. *Virus Evolution* **7**. ISSN 2057-1577. (doi:10.1093/ve/veab064).  
360 Veab064.
- 361 [20] Katoh, K. & Standley, D. M., 2013 MAFFT multiple sequence alignment software version 7:  
362 improvements in performance and usability. *Molecular biology and evolution* **30**, 772–780. ISSN  
363 0737-4038. (doi:10.1093/molbev/mst010).
- 364 [21] Minh, B. Q., Schmidt, H. A., Chernomor, O., Schrempf, D., Woodhams, M. D., von Haeseler, A.  
365 & Lanfear, R., 2020 Iq-tree 2: New models and efficient methods for phylogenetic inference in the  
366 genomic era. *Molecular Biology and Evolution* **37**, 1530–1534. (doi:10.1093/molbev/msaa015).
- 367 [22] Sagulenko, P., Puller, V. & Neher, R. A., 2018 TreeTime: Maximum-likelihood phylogenetic  
368 analysis. *Virus Evolution* **4**. ISSN 2057-1577. (doi:10.1093/ve/vex042). Vex042.

- 369 [23] Rambaut, A., 2020. Phylodynamic Analysis | 176 genomes | 6 Mar 2020. Library Catalog:  
370 virological.org.
- 371 [24] Zhou, F., Yu, T., Du, R., Fan, G., Liu, Y., Liu, Z., Xiang, J., Wang, Y., Song, B., Gu, X. *et al.*, 2020  
372 Clinical course and risk factors for mortality of adult inpatients with COVID-19 in Wuhan, China:  
373 a retrospective cohort study. *The Lancet* ISSN 0140-6736. (doi:10.1016/S0140-6736(20)30566-3).
- 374 [25] Rambaut, A., Drummond, A. J., Xie, D., Baele, G. & Suchard, M. A., 2018 Posterior Summariza-  
375 tion in Bayesian Phylogenetics Using Tracer 1.7. *Systematic Biology* **67**, 901–904. ISSN 1063-5157.  
376 (doi:10.1093/sysbio/syy032).
- 377 [26] Cori, A., Ferguson, N. M., Fraser, C. & Cauchemez, S., 2013 A New Framework and Software to  
378 Estimate Time-Varying Reproduction Numbers During Epidemics. *Am J Epidemiol* **178**, 1505–  
379 1512. ISSN 0002-9262. (doi:10.1093/aje/kwt133).
- 380 [27] Thompson, R. N., Stockwin, J. E., van Gaalen, R. D., Polonsky, J. A., Kamvar, Z. N., De-  
381 marsh, P. A., Dahlgvist, E., Li, S., Miguel, E., Jombart, T. *et al.*, 2019 Improved inference of  
382 time-varying reproduction numbers during infectious disease outbreaks. *Epidemics* **29**, 100356.  
383 (doi:10.1016/j.epidem.2019.100356).
- 384 [28] Alizon, S., Haim-Boukobza, S., Foulongne, V., Verdurme, L., Trombert-Paolantoni, S., Lecorche,  
385 E., Roquebert, B. & Sofonea, M. T., 2021 Rapid spread of the SARS-CoV-2 Delta vari-  
386 ant in some French regions, June 2021. *Eurosurveillance* **26**, 2100573. (doi:10.2807/1560-  
387 7917.ES.2021.26.28.2100573).
- 388 [29] Sofonea, M. T., Roquebert, B., Foulongne, V., Morquin, D., Verdurme, L., Trombert-Paolantoni,  
389 S., Roussel, M., Bonetti, J.-C., Zerah, J., Haim-Boukobza, S. *et al.*, 2022 Analyzing and Modeling  
390 the Spread of SARS-CoV-2 Omicron Lineages BA.1 and BA.2, France, September 2021–February  
391 2022. *Emerging Infectious Diseases* **28**. (doi:10.3201/eid2807.220033).
- 392 [30] Martin, M. A., VanInsberghe, D. & Koelle, K., 2021 Insights from SARS-CoV-2 sequences. *Science*  
393 **371**, 466–467. (doi:10.1126/science.abf3995).

- 394 [31] Wilkinson, E., Giovanetti, M., Tegally, H., San, J. E., Lessells, R. & *et al.*, 2021 A year of genomic  
395 surveillance reveals how the SARS-CoV-2 pandemic unfolded in Africa. *Science* **374**, 423–431.  
396 (doi:10.1126/science.abj4336).
- 397 [32] Suchard, M. A., Lemey, P., Baele, G., Ayres, D. L., Drummond, A. J. & Rambaut, A., 2018  
398 Bayesian phylogenetic and phylodynamic data integration using BEAST 1.10. *Virus Evolution* **4**.  
399 (doi:10.1093/ve/vey016).
- 400 [33] Hall, M. D., Woolhouse, M. E. J. & Rambaut, A., 2016 The effects of sampling strategy on the  
401 quality of reconstruction of viral population dynamics using Bayesian skyline family coalescent  
402 methods: A simulation study. *Virus Evolution* **2**, vew003. (doi:10.1093/ve/vew003).
- 403 [34] Karcher, M. D., Carvalho, L. M., Suchard, M. A., Dudas, G. & Minin, V. N., 2020 Estimating  
404 effective population size changes from preferentially sampled genetic sequences. *PLOS Computa-*  
405 *tional Biology* **16**, e1007774. (doi:10.1371/journal.pcbi.1007774).
- 406 [35] Guindon, S. & De Maio, N., 2021 Accounting for spatial sampling patterns in Bayesian  
407 phylogeography. *Proceedings of the National Academy of Sciences* **118**, e2105273118.  
408 (doi:10.1073/pnas.2105273118).
- 409 [36] Louca, S., McLaughlin, A., MacPherson, A., Joy, J. B. & Pennell, M. W., 2021 Fundamental  
410 Identifiability Limits in Molecular Epidemiology. *Molecular Biology and Evolution* **38**, 4010–4024.  
411 (doi:10.1093/molbev/msab149).