

## Using multiple sampling strategies to estimate SARS-CoV-2 epidemiological parameters from genomic sequencing data

### AUTHOR LIST AND AFFILIATIONS

Rhys P. D. Inward<sup>1,6</sup>, Kris V. Parag<sup>2,3,5,6</sup>, Nuno R. Faria<sup>1,2,4,5,6</sup>

1. Department of Zoology, University of Oxford, Oxford, UK
2. MRC Centre of Global Infectious Disease Analysis, Jameel Institute for Disease and Emergency Analytics, Imperial College London, London, UK
3. NIHR Health Protection Research Unit in Behavioural Science and Evaluation, University of Bristol, Bristol, UK
4. Instituto de Medicina Tropical, Faculdade de Medicina da Universidade de Sao Paulo, Sao Paulo, Brazil
5. Jointly supervised this work
6. Corresponding author: E-mail: [rhys.inward@zoo.ox.ac.uk](mailto:rhys.inward@zoo.ox.ac.uk), [k.parag@imperial.ac.uk](mailto:k.parag@imperial.ac.uk), [n.faria@imperial.ac.uk](mailto:n.faria@imperial.ac.uk)

**NOTE: This preprint reports new research that has not been certified by peer review and should not be used to guide clinical practice.**

## ABSTRACT

SARS-CoV-2 virus genomes are currently being sequenced at an unprecedented pace. The choice of sequences used in genetic and epidemiological analysis is important as it can induce biases that detract from the value of these rich datasets. This raises questions about how a set of sequences should be chosen for analysis, and which epidemiological parameters derived from genomic data are sensitive or robust to changes in sampling. We provide initial insights on these largely understudied problems using SARS-CoV-2 genomic sequences from Hong Kong and the Amazonas State, Brazil. We consider sampling schemes that select sequences uniformly, in proportion or reciprocally with case incidence and which simply use all available sequences (unsampled). We apply *Birth-Death Skyline* and *Skygrowth* methods to estimate the time-varying reproduction number ( $R_t$ ) and growth rate ( $r_t$ ) under these strategies as well as related  $R_0$  and date of origin parameters. We compare these to estimates from case data derived from *EpiFilter*, which we use as a reference for assessing bias. We find that both  $R_t$  and  $r_t$  are sensitive to changes in sampling whilst  $R_0$  and date of origin are relatively robust. Moreover, we find that the unsampled datasets (opportunistic sampling) provided, overall, the worst  $R_t$  and  $r_t$  estimates for both Hong Kong and the Amazonas case studies. We highlight that sampling strategy may be an influential yet neglected component of sequencing analysis pipelines. More targeted attempts at genomic surveillance and epidemic analyses, particularly in resource-poor settings which have a limited genomic capability, are necessary to maximise the informativeness of virus genomic datasets.

## INTRODUCTION

Severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2) is an enveloped single-stranded zoonotic RNA virus belonging to the *Betacoronavirus* genus and *Coronaviridae* family (Gorbalenya *et al.*, 2020). It was first identified in late 2019 in a live food market in Wuhan City, Hubei Province, China (Zhu *et al.*, 2020). Within a month, SARS-CoV-2 had disseminated globally through sustained human-to-human transmission. It was declared a public health emergency of international concern on the 30th of January 2020 by the World Health Organisation (World Health Organisation, 2020). Those infected with SARS-CoV-2 have phenotypically diverse symptoms ranging from mild fever to multiple organ dysfunction syndromes and death (Verity *et al.*, 2020).

Despite the implementation of non-pharmaceutical interventions (NPIs) by many countries to control their epidemics, to date over 300 million SARS-CoV-2 cases and 5.4 million deaths have been reported worldwide (World Health Organisation, 2022). These NPIs can vary within and between countries and include restrictions on international and local travel, school closures, social distancing measures and the isolation of infected individuals and their contacts (European Centre for Disease Prevention and Control, 2020). The key aim of NPIs is to reduce epidemic transmission, often measured by epidemiological parameters such as the time-varying reproduction number ( $R_t$  at time  $t$ ) and growth rate ( $r_t$ ) (Supplementary Table 1) (Anderson *et al.*, 2020; UK Health Security Agency, 2022). However, there is currently great difficulty in estimating and comparing epidemiological parameters derived from case and death data globally due to disparities in molecular diagnostic surveillance and notification systems between countries. Further, even if data are directly comparable, the choice of epidemiological parameter can implicitly shape insights into how NPIs influence transmission potential (Dushoff and Park, 2021; Parag, Thompson and Donnelly, 2021). As such, there is a need to use alternative data sources, such as genomic data (World Health Organisation, 2021a), to gain improved insights into viral transmission dynamics (Jombart *et al.*, 2014; Duchene *et al.*, 2020).

Phylogenetic analysis of virus genome sequences have increasingly been used for studying emerging infectious diseases, as seen during the current SARS-CoV-2 pandemic (Faria *et al.*, 2021; Nadeau *et al.*, 2021; Romano and Melo, 2021; Volz *et al.*, 2021), recent Ebola virus outbreaks in Western Africa (Dudas *et al.*, 2017) and Zika outbreaks in Brazil and the

Americas (Faria *et al.*, 2017; Grubaugh *et al.*, 2017). Transmissibility estimates such as the basic reproduction number ( $R_0$ ),  $R_t$  and  $r_t$  can be directly inferred from genomic sequencing data in addition to other epidemiological parameters like the date of origin of a given viral variant which can only be inferred from genomic data. This is of particular importance for variants of concern (VOC), genetic variants with evidence of increased transmissibility, more severe disease, and/ or immune evasion. VOC are typically detected through virus genome sequencing and there is often a limited understanding of their epidemiological characteristics from epidemiological data alone (Harvey *et al.*, 2021). To maximise the use of additional epidemiological information from genomic data, clear guidelines on sampling need to be provided (Lesley *et al.*, 2021).

Currently, SARS-CoV-2 virus genomes from COVID-19 cases are being sequenced at an unprecedented pace providing a wealth of virus genomic datasets (Rambaut *et al.*, 2020). There are currently over 7.4 million genomic sequences available on GISAID, an open-source repository for influenza and SARS-CoV-2 genomic sequences (Shu and McCauley, 2017). These rich datasets can be used to provide an independent perspective and can help validate or challenge parameters derived from epidemiological data. Moreover, the use of genomic data can overcome some of the limitations and biases of using epidemiological data alone. For example, it is less susceptible to changes at the government level such as alterations to the definition of a confirmed case and changes to notification systems (de Souza *et al.*, 2020; Tsang *et al.*, 2020). Inferences from virus genomic data improve our understanding of underlying epidemic spread and can facilitate better-informed infection control decisions (Dolan, Whitfield and Andino, 2018).

The most popular approaches used to investigate changes in virus population dynamics include the Bayesian Skyline Plot (Drummond *et al.*, 2005) and Skygrid (Gill *et al.*, 2013) models and the birth-death skyline (BDSKY) (Stadler *et al.*, 2013). These integrate Markov Chain Monte Carlo (MCMC) procedures and often converge slowly on large datasets (Hall, Woolhouse and Rambaut, 2016). As such, currently available SARS-CoV-2 datasets containing thousands of sequences become computationally impractical to analyse and sub-sampling is necessary. Although there have been some previous studies (Stack *et al.*, 2010; de Silva, Ferguson and Fraser, 2012; Hall, Woolhouse and Rambaut, 2016; Karcher *et al.*, 2016; Parag, du Plessis and Pybus, 2020), the effects of sampling strategies on phylogenetic and phylodynamic inferences of pathogens is currently a neglected area of study (Frost *et al.*,

2015), particularly concerning SARS-CoV-2. To our knowledge, there are no published studies concerning SARS-CoV-2 which explore the effect that sampling strategies have on the phylodynamic reconstruction of key transmission parameters. This is important as incorrectly implementing a sampling scheme or ignoring its importance can mislead inferences and introduce biases (Hall, Woolhouse and Rambaut, 2016; Hidano and Gates, 2019). This raises the important question of how a set of sequences should be selected for analysis and which parameters are sensitive or robust to changes in sampling.

Here we aim to explore how diverse sampling strategies in genomic sequencing may affect the estimation of key epidemiological parameters from genomic data. To do this, we estimate  $R_0$ ,  $R_t$ , and  $r_t$  from genomic sequencing data under different sampling strategies from a location with high genomic coverage represented by Hong Kong, and a location with low genomic coverage represented by the Amazonas region, Brazil. Moreover, we compare epidemiological parameters derived from genomic data to those estimated from corresponding epidemiological data which we considered here as our gold standard. By getting genomic inferences close to the case data we can then draw better inferences of transmission estimates and parameters that cannot be derived from case data alone. This will help us to understand the impact that sampling strategies have on phylodynamic inference and aid in the interpretation of epidemiological parameters from areas with differing genomic coverage.

## METHODS

### Empirical Estimation of the Reproduction Number, Time-varying Effective Reproduction Number, and Growth Rate

#### *Epidemiological Datasets*

Two sources of data from the Amazonas region, Brazil and one source of data from Hong Kong were used in calculating empirical epidemiological parameters. For the Amazonas region, mortality, and case data from the SIVEP-Gripe (Sistema de Informação de Vigilância Epidemiológica da Gripe) SARI (severe acute respiratory infections) database, including both class 4 and 5 death records (corresponding to confirmed and suspected COVID-19 deaths), from the 30<sup>th</sup> of November 2020 up to 7<sup>th</sup> of February 2021, were used. Here we were interested in cases caused by the P.1/Gamma VOC first detected in Manaus, the number of P.1 cases was calculated by using the proportion of P.1/Gamma viral sequences uploaded to GISAID within each week (Supplementary Figure 1). For Hong Kong, all case and mortality data were extracted from the Centre of Health Protection, Department of Health, the Government of the Hong Kong Special Administrative region up to the 7<sup>th</sup> of May 2020. Due to lags in the development of detectable viral loads, symptom onset and subsequent testing (Gostic *et al.*, 2020); the date in which each case was recorded was left shifted by 5 days within our models (Pullano *et al.*, 2021) to account for these delays in both datasets.

#### *Basic Reproduction Number*

The  $R_0$  was estimated using a time series of confirmed SARS-CoV-2 cases from both Hong Kong and the Amazonas region. To avoid the impact of NPIs on  $R_0$  estimates, only data up to the banning of mass gathering in Hong Kong (27<sup>th</sup> March 2020) and up to the imposition of strict restrictions in the Amazonas region (12<sup>th</sup> January 2021) were used. Weekly counts of confirmed cases were modelled using maximum likelihood methods. The weekly case counts were assumed to be Poisson distributed and were fitted to a deterministic closed Susceptible-Exposed-Infectious-Recovered (SEIR) model (Equation 1) by maximising the likelihood of observing the data given the model parameters (Table 1).

Equation 1:

$$\lambda = \frac{\beta(I)}{N} \frac{dS}{dt} = -\lambda S \frac{dE}{dt} = \lambda S - \gamma E \frac{dI}{dt} = \gamma E - \sigma I \frac{dR}{dt} = \sigma I$$

Subsequently, the log-likelihood was used to calculate the  $R_0$  by fitting  $\beta$ , the effective contact rate (Equation 2).

Equation 2:

$$R_0 = \beta\alpha$$

To generate approximate confidence intervals for  $R_0$ , bootstrapping was used with 1000 iterations.

**Table 1:** This shows the parameter estimates used within the deterministic SEIR model.

Parameter	Description	Value (source)
$R_0$	Basic Reproduction Number	Estimated
N	Population of Hong Kong	7,481,800 persons (The World Bank, 2021)
	Population of Amazonas Region	4,207,714 persons (IBGE, 2020)
$\beta$	Effective Contact Rate	Estimated
$\alpha$	Infectious Period	0.07 (Byrne <i>et al.</i> , 2020)
$\lambda$	Force of Infection	Estimated
$\gamma$	Progression from E to I	5.26 day <sup>-1</sup> (McAloon <i>et al.</i> , 2020)
$\delta$	Progression from I to R	14.3 day <sup>-1</sup> (Byrne <i>et al.</i> , 2020)
S	Susceptible compartment	Estimated

E	Exposed Compartment	Estimated
I	Infectious Compartment	Estimated
R	Recovered Compartment	Estimated

### *Time-varying Effective Reproduction Number*

To estimate the  $R_t$  from empirical line list data the *EpiFilter* model (Parag, 2021) was used. To estimate  $R_t$ , *EpiFilter* uses a renewal transmission model; a general and popular framework used in the modelling of infectious diseases (Fraser, 2007). This model describes how the number of new cases (incidence) at time  $t$  depends on the  $R_t$  at that specified time point and the past incidence, which is summarised by the cumulative number of cases up to each time point weighted by the generation time distribution. Moreover, *EpiFilter* integrates both Bayesian forward and backward recursive smoothing. This improves  $R_t$  estimates by leveraging the benefits of two of the most popular  $R_t$  estimation approaches: *EpiEstim* (Cori *et al.*, 2013) and the Wallinga-Teunis equation (Wallinga and Teunis, 2004). Both methods only utilise a proportion of the information available with either past or future incidence being informative. *EpiFilter* combines both past and future information and consequently minimises the mean squared error in estimation and reduces dependence on prior assumptions. We assume the generation time distribution is well approximated by the serial interval (SI) distribution (Flaxman *et al.*, 2020). *EpiFilter* was used as a reference for parameters estimated from genomic data.

### *Growth Rate*

After the  $R_t$  has been inferred, its relationship with  $r_t$  as described by the Wallinga-Lipsitch equation for a gamma distributed generation time (Equation 3) was used to estimate  $r_t$  (Wallinga and Lipsitch, 2007). The SI and variance for Hong Kong were derived from a systematic review and meta-analysis exploring these values (Rai, Shukla and Dwivedi, 2021) and a study exploring SI in Brazil was used for the Amazonas datasets (Prete *et al.*, 2021). The SI was assumed to be gamma distributed. The gamma distribution is represented by  $\text{gamma} = (\epsilon, \gamma)$ .



Equation 3:

$$r_t = \varepsilon(R_t^{\left(\frac{1}{\gamma}\right)} - 1)$$

### **SARS-CoV-2 Brazilian Gamma VOC and Hong Kong datasets**

All high-quality, complete SARS-CoV-2 genomes were downloaded from GISAID (Shu and McCauley, 2017) for Hong Kong (up to 7<sup>th</sup> May 2020) and the Amazonas state, Brazil (from 30<sup>th</sup> November 2020 up to 7<sup>th</sup> February 2021). Using the Accession ID of each sequence, all sequences were screened and only sequences previously analysed and published in PubMed, MedRxiv, BioRxiv, virological or Preprint repositories were selected for subsequent analysis. For both datasets, sequence alignment was conducted using MAFFT.V.7 (Kato *et al.*, 2002). The first 130 base pairs (bp) and last 50 bps of the aligned sequences were trimmed to remove potential sequencing artefacts in line with the Nextstrain protocol (Hadfield *et al.*, 2018). Both datasets were then processed using the Nextclade pipeline for quality control (<https://clades.nextstrain.org/>). Briefly, the Nextclade pipeline examines the completeness, divergence, and ambiguity of bases in each genetic sequence. Only sequences deemed ‘good’ by the Nextclade pipeline were selected for. Subsequently, all sequences were screened for identity and in the case of identical sequences, for those with the same location, collection date, only one such isolate was used. Moreover, PANGO lineage classification was conducted using the Pangolin (Rambaut *et al.*, 2020) software tool (<http://pangolin.cog-uk.io>) on sequences from the Amazonas region and only those with the designated P.1/Gamma lineage were selected for (Supplementary Figure 1).

### **Maximum Likelihood tree reconstruction**

Maximum likelihood phylogenetic trees were reconstructed using IQTREE2 (Minh *et al.*, 2020) for both datasets. A TIM2 model of nucleotide substitution with empirical base frequencies and a proportion of invariant sites was used as selected for by the ModelFinder application (Kalyaanamoorthy *et al.*, 2017) for the Hong Kong dataset. For the Brazilian dataset, a TN model of nucleotide substitution (Tamura and Nei, 1993) with empirical base frequencies was selected for. To assess branch support, the approximate likelihood-ratio test based on the Shimodaira–Hasegawa-like procedure with 1,000 replicates (Anisimova *et al.*, 2011), was used.

### Root-to-tip regression

To explore the temporal structure of both the Brazilian and Hong Kong dataset, TempEst v.1.5.3 (Rambaut *et al.*, 2016) was used to regress the root-to-tip genetic distances against sampling dates (yyyy-mm-dd). The ‘best-fitting’ root for the phylogeny was found by maximising the  $R^2$  value of the root-to-tip regression. Several sequences showed incongruent genetic diversity and were discarded from subsequent analyses. This resulted in a final dataset of  $N = 117$  Hong Kong sequences and  $N = 196$  Brazilian sequences. The gradient of the slopes (clock rates) provided by TempEst were used to inform the clock prior in the phylodynamic analysis.

### Subsampling for analysis

Four retrospective sampling schemes were used to select a subsample of Amazonas and Hong Kong sequences. Each sampling period was broken up into weeks with each week being used as an interval according to a temporal sampling scheme (without replacement). This temporal sampling scheme was based on the number of reported cases of SARS-CoV-2.

Temporal sampling schemes explored were:

- **Uniform sampling:** All weeks have equal probability.
- **Proportional sampling:** Weeks are chosen with a probability proportional to the value of the number of cases in each epi-week.
- **Reciprocal-proportional sampling:** Weeks are chosen with a probability proportional to the reciprocal of the number of cases in each epi-week.
- **No sampling strategy applied:** All sequences were included without a sampling strategy applied.

These sampling schemes were inspired by those recommended by the WHO for practical use in different settings and scenarios (World Health Organisation, 2021b). Proportional sampling is equivalent to representative sampling, uniform sampling is equivalent to fixed sampling whilst the unsampled data includes all sampling strategies. Reciprocal-proportional sampling is not commonly used in practice as was used as a control within this study.

### Bayesian Evolutionary Analysis

Date molecular clock phylogenies were inferred for all sampling strategies applied to the Amazonas and Hong Kong dataset using BEAST v1.10.4 (Suchard *et al.*, 2018) with

BEAGLE library v3.1.0 (Ayres *et al.*, 2019) for accelerated likelihood evaluation. For both the Amazonas and Hong Kong datasets, a HKY substitution model with gamma-distributed rate variation among sites and four rate categories was used to account for among-site rate variation (Hasegawa, Kishino and Yano, 1985). A strict clock molecular clock model was chosen. Both the Amazonas and Hong Kong dataset were analysed under a flexible non-parametric skygrid tree prior (Hill and Baele, 2019). Four independent MCMC chains were run for both datasets. For the Amazonas dataset, each MCMC chain consisted of 250,000,000 steps with sampling every 50,000 steps. Meanwhile, for the Hong Kong dataset, each MCMC chain consisted of 200,000,000 steps with sampling every 40,000 steps. For both datasets, the four independent MCMC runs were combined using LogCombiner v1.10.4 (Suchard *et al.*, 2018). Subsequently, 10% of all trees were discarded as burn in, and the effective sample size of parameter estimates were evaluated using TRACER v1.7.2 (Rambaut *et al.*, 2018). An effective sample size of over 200 was obtained for all parameters. Maximum clade credibility (MCC) trees were summarised using Tree Annotator (Suchard *et al.*, 2018).

## **Phylodynamic Reconstruction**

### *Estimation of the Reproduction Number and Time-varying Effective Reproduction Number*

The Bayesian birth-death skyline (BDSKY) model (Stadler *et al.*, 2013) implemented within BEAST 2 v2.6.5 (Bouckaert *et al.*, 2019) was used to estimate time-varying rates of epidemic transmission, measured as changes in  $R_t$  (Table 2). A HKY substitution model with a gamma-distributed rate variation among sites and four rate categories (Hasegawa, Kishino and Yano, 1985) was used alongside a strict molecular clock model. A lognormal distribution was used for  $R_t$ . The selected number of intervals for both datasets was 5, representing  $R_t$  changing every 2.5 weeks for the Hong Kong datasets and every 2 weeks for the Brazilian datasets, with equidistant intervals per step. An exponential distribution was used with a mean of  $36.5y^{-1}$  for the rate of becoming infectious, assuming a mean duration of infection of 10 days (Nadeau *et al.*, 2021). A uniform distribution was used for the sampling proportion. Four independent MCMC chains were run for 50 million MCMC steps with sampling every 5000 steps for each dataset. The four independent MCMC runs were combined using LogCombiner v2.6.5. (Bouckaert *et al.*, 2019) and the effective sample size of parameter estimates were evaluated using TRACER v1.7.2 (Rambaut *et al.*, 2018). An effective sample size of over 200 was obtained for all parameters. The *bdskytools* R package (<https://github.com/laduplessis/bdskytools>) was used to plot the BDSKY results.

**Table 2:** Values and priors for the parameters of the BDSKY model

Parameter	Dataset	Value or prior	Rationale/Assumption
Clock rate	Brazil	$4.0 \times 10^{-4}$ (subs/site/year)	Informed by root-to-tip regression
	Hong Kong	$1.0 \times 10^{-4}$ (subs/site/year)	
Death rate	Brazil and Hong Kong	$36.5 \text{ y}^{-1}$	The period between infection and becoming uninfected assumed an exponential distribution with a mean of 10 d (Nadeau <i>et al.</i> , 2021)
Reproductive number	Brazil and Hong Kong	Lognormal (0.8, 0.5)	Median 2.2, 95% IQR 0.8 to 5.9
Time of origin	Brazil	Lognormal (-1.50, 0.4) y before present	Median 4 <sup>th</sup> December 2020, 95% IQR 25 <sup>th</sup> September 2020 to 12 <sup>th</sup> January, 2021
	Hong Kong	Lognormal (-1.75, 0.4) y before present	Median 18 <sup>th</sup> January 2020, 95% IQR 17 <sup>th</sup> November 2019 to 15 <sup>th</sup> February 2020
Sampling proportion	Brazil	Uniform (0, 0.024)	196 sequences from 8246 suspected P.1 cases as of 7 <sup>th</sup> February, 2021

	Hong Kong	Uniform (0, 0.116)	117 sequences from 1012 confirmed cases as of 7 <sup>th</sup> May, 2020
--	-----------	--------------------	---

### *Estimation of Growth Rates*

For each dataset, a scaled proxy for  $r_t$  was estimated through time using the *skygrowth* model (Volz and Didelot, 2018) within R. *Skygrowth* uses MCMC to apply a first-order autoregressive stochastic process, founded on a non-parametric Bayesian approach, on the growth rate of the effective population size. The MCMC chains were run for one million iterations for each dataset on their MCC tree with an Exponential ( $10^{-5}$ ) prior on the smoothing parameter. The *skygrowth* model was parameterised assuming that the effective population size of SARS-COV-2 could change every two weeks. To enable comparisons of  $r_t$  estimated by *skygrowth* and  $r_t$  estimated by *EpiFilter*, the  $r_t$  provided by the *skygrowth* model was converted to the exponential growth rate. To do this, the  $R_t$  was calculated from  $r_t$  by adding a gamma rate variable which assumed a mean duration of infection of 10 days (Nadeau *et al.*, 2021). Subsequently, the Wallinga-Lipsitch equation (Equation 3) was used to convert  $R_t$  into the exponential growth rate (Wallinga and Lipsitch, 2007).

### **Comparing Parameters Estimates from Genetic and Epidemiological Data**

To compare parameters estimates from epidemiological and genetic data the Jensen-Shannon divergence ( $D_{JS}$ ) (Lin, 1991), which measures the similarity between two probability mass functions (PMFs), was applied. The  $D_{JS}$  offers a formal information theoretic evaluation of distributions and is more robust than comparing Bayesian credible intervals (BCIs) since it considers both the shape and spread of a given distribution. The  $D_{JS}$  is essentially a symmetric and smoothed version of the Kullback-Leibler divergence ( $D_{KL}$ ) and is commonly used in the fields of machine learning and bioinformatics. The  $D_{KL}$  between two PMFs, P and Q, is defined as  $D_{KL}$  in Equation 4 below (Kullback and Leibler, 1951).

Equation 4:

$$D_{KL} (P || M) = \sum_{x \in X} P(x) \log\left(\frac{P(x)}{Q(x)}\right)$$

To calculate the PMF for each epidemiological parameter, the cumulative probability density function (PDF) was extracted for each model, converted to a probability density function (PDF), and a discretisation procedure then applied (Equation 5).  $\tau$  represents the PDF and is discretized via Equation 4, where  $s = 0.05, 0.01, \dots$  and  $\tau(v)$  is the cumulative probability density of  $\tau$ .

Equation 5:

$$\tau_{Rt,rt,R0} = \int_{s-0.025}^{s+0.025} \tau(v)$$

The Jensen-Shannon distance (JSD) metric quantifies the square-root of the total  $D_{JS}$  to the average probability distribution and is the metric that we used to compare parameter estimations from differing sampling strategies. The  $D_{JS}$  can be calculated using Equation 6 with P and Q representing the two probability distributions and  $D_{KL}$  as the KL divergence. A smaller JSD metric indicates that P and Q are more similar with a Jensen-Shannon distance of 0 indicating equivalence of the two distributions. The mean JSD was taken over all intervals for the BDSKY and *Skygrowth* models to obtain an overall measure of the level of estimated similarity.

Equation 6:

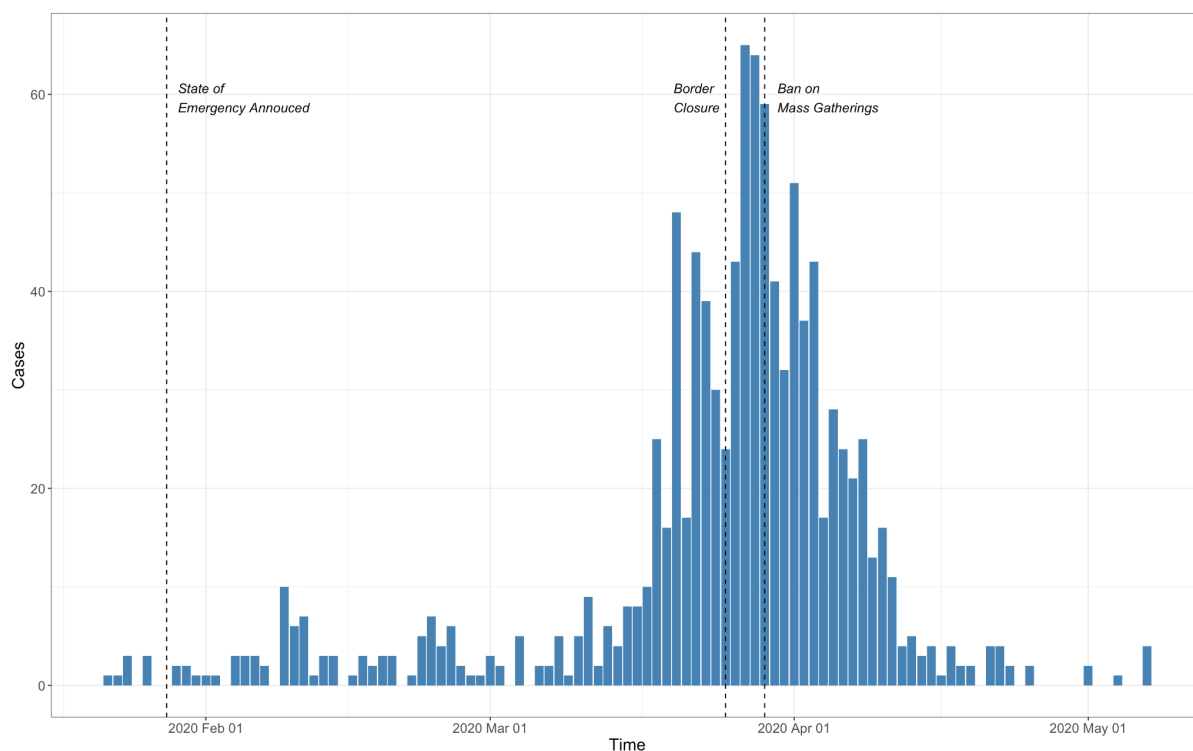
$$D_{JS}(P \parallel Q) = \frac{1}{2}D_{KL}(P \parallel M) + \frac{1}{2}D_{KL}(Q \parallel M) \text{ where } M = \frac{1}{2}(P + Q)$$

## RESULTS

### Sampling Schemes

#### *Hong Kong*

Hong Kong reacted rapidly upon learning of the emergence of SARS-CoV-2 in Wuhan, Hubei province, China by declaring a state of emergency on the 25th of January 2020 and by mobilising intensive surveillance schemes in response to initial cases (Cowling *et al.*, 2020). This appeared to be successful in controlling the first wave of cases. However, due to imported cases from Europe and North America, a second wave of SARS-CoV-2 infections emerged prompting stricter NPIs such as the closure of borders and restrictions on gatherings (Cowling *et al.*, 2020). Following these measures, the incidence of SARS-CoV-2 rapidly decreased (Figure 1). Hong Kong has a high sampling intensity with 11.6% of confirmed cases sequenced during our study period.



**Figure 1:** Confirmed SARS-CoV-2 cases from Hong Kong until 7<sup>th</sup> of May 2020. The dashed lines represent policy change-times (Cowling *et al.*, 2020).

The number of cases within Hong Kong for each week was used to inform the sampling schemes used within this study. This resulted in the unsampled scheme having  $N = 117$  sequences, the proportional sampling scheme having  $N = 54$  sequences, the uniform sampling

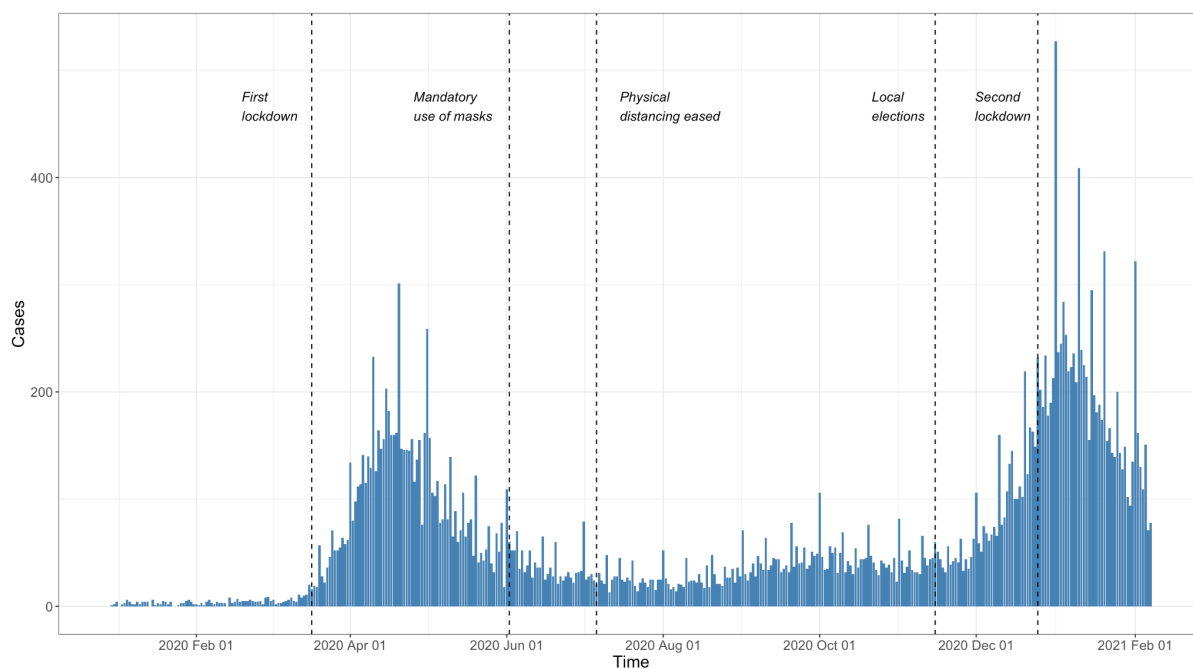
scheme having  $N = 79$  and the reciprocal-proportional sampling scheme having  $N = 84$  sequences (Supplementary Figure 2).

### *Amazonas*

The Amazonas state of Brazil had its first laboratory confirmed case of SARS-CoV-2 in March 2020 in a traveller returning from Europe (Nascimento *et al.*, 2020). The first wave of SARS-CoV-2 infections within the state peaked in early May 2020 (Figure 2). From then, the epidemic waned, cases dropped, remaining stable until mid-December 2020. The number of cases then started growing exponentially, ushering in a second epidemic wave. This second wave peaked in January 2021 (Figure 2) and was caused by the emergence of a new SARS-CoV-2 VOC, designated P.1/Gamma (Faria *et al.*, 2021).

To combat this second wave, the Government of the Amazonas state suspended all non-essential commercial activities on the 23rd of December 2020 (<http://www.pge.am.gov.br/legislacao-covid-19/>). However, in response to protests, these restrictions were reversed, and cases continued to climb. On the 12th of January, NPIs were re-introduced (<http://www.pge.am.gov.br/legislacao-covid-19/>) which seemed to be successful in reducing the case incidence in the state. However, cases remain comparatively high (Figure 2). Amazonas has a low sampling intensity with 2.4% of suspected P.1/gamma cases sequenced during our study period.





**Figure 2:** Confirmed SARS-CoV-2 cases from Amazonas state, north Brazil until 7<sup>th</sup> of February 2021. The dashed lines represent policy change-times (Sabino *et al.*, 2021).

The number of cases within the Amazonas region informed the sampling schemes used within this study. This resulted in the unsampled scheme having  $N = 196$  sequences, the proportional sampling scheme having  $N = 168$  sequences, the uniform sampling scheme having  $N = 150$  and the reciprocal-proportional sampling scheme having  $N = 67$  sequences (Supplementary Figure 3).

### Root-to-tip Regression

The correlation ( $R^2$ ) between genetic divergence and sampling dates for the Hong Kong datasets ranged between 0.36 and 0.52 and between 0.13 and 0.20 for the Amazonas datasets. This implies that the Hong Kong datasets have a stronger temporal signal. This is likely due to the Hong Kong datasets have a wider sampling interval (106 days) compared to the Amazonas datasets (69 days). A wider sampling interval can lead to a stronger temporal signal (Drummond *et al.*, 2003). No association between the number of sequences in each sampling scheme and the  $R^2$  was found. This implies that the data has a high degree of non-independence which is an unexpected finding as more independent data should reduce the effects of stochasticity. The gradient (rate) of the regression ranged from  $1.24 \times 10^{-3}$  to  $1.72 \times 10^{-3}$  s/s/y for the Hong Kong datasets and  $4.41 \times 10^{-4}$  to  $5.28 \times 10^{-4}$  s/s/y for the Amazonas datasets.

## Estimation of Evolutionary Parameters

The mean substitution rate (measured in units of number of substitutions per site per year,  $s/s/y$ ) and the time to most common recent ancestor (TMRCA) was estimated in BEAST, for both datasets, and the estimation from all sampling schemes was compared.

### *Hong Kong*

For Hong Kong, the mean substitution rate per site per year ranged from  $9.16 \times 10^{-4}$  to  $2.09 \times 10^{-3}$  with sampling schemes all having overlapped BCI (Supplementary table 2; Supplementary Figure 4A). This indicates that the sampling scheme did not have a significant impact on the estimation of the clock rate. Moreover, the clock rate is comparable to estimations from the root-to-tip regression and to early estimations of the mean substitution rate per site per year of SARS-CoV-2 (Duchene et al., 2020).

Molecular clock dating of the Hong Kong dataset indicates that the estimated time of the most common recent ancestor was mid-November 2019 and early January 2020 (mean, 10th December 2019; 95% BCI interval, 14th November 2019 – 1st January 2020, Figure 3B; Supplementary Table 2). This is around 5 weeks before the first confirmed case which was reported on the 18th of January 2021. Once again, all sampling strategies have overlapped BCIs suggesting that the sampling scheme does not significantly impact the estimation of the TMRCA.

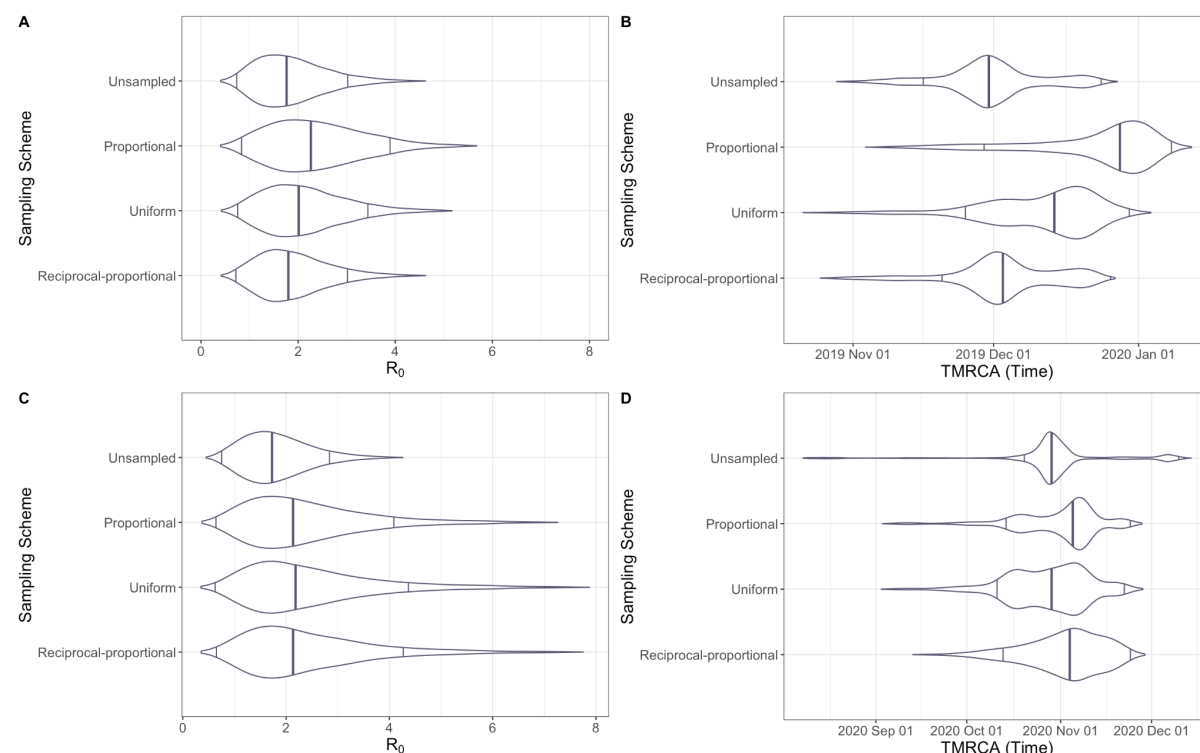
### *Brazil*

For the P.1 lineage in the Amazonas region, the mean substitution rate ranged from  $4.00 \times 10^{-4}$  to  $5.56 \times 10^{-4}$  with all sampling schemes having overlapped BCIs (Figure 3D, Supplementary Table 2; Supplementary Figure 4B). This indicates that sampling strategy does not impact the estimation of the clock rate, supporting findings from the Hong Kong dataset. This supports estimations from the root-to-tip analysis.

Molecular clock dating estimated a TMRCA between mid-September and mid-November (mean, 23rd October 2020; 95% BCI interval, 16th September 2020 – 18th November 2020, Figure 3D; Supplementary Table 2). This is around five weeks before the date of the first P.1 case identified in Manaus used in our study. All sampling schemes have overlapping BCI supporting the inference from the Hong Kong datasets that TMRCA is robust to sampling.

### Estimation of Basic Reproduction Number

We found that Hong Kong had a significantly lower  $R_0$  of 2.17 (95% credible interval (CI) = 1.43 - 2.83) when compared to Amazonas which had a  $R_0$  of 3.67 (95% CI = 2.83 – 4.48). All sampling schemes for both datasets were characterised by similar  $R_0$  values (Figure 3) indicating that the estimation of  $R_0$  is robust to changes in sampling scheme.



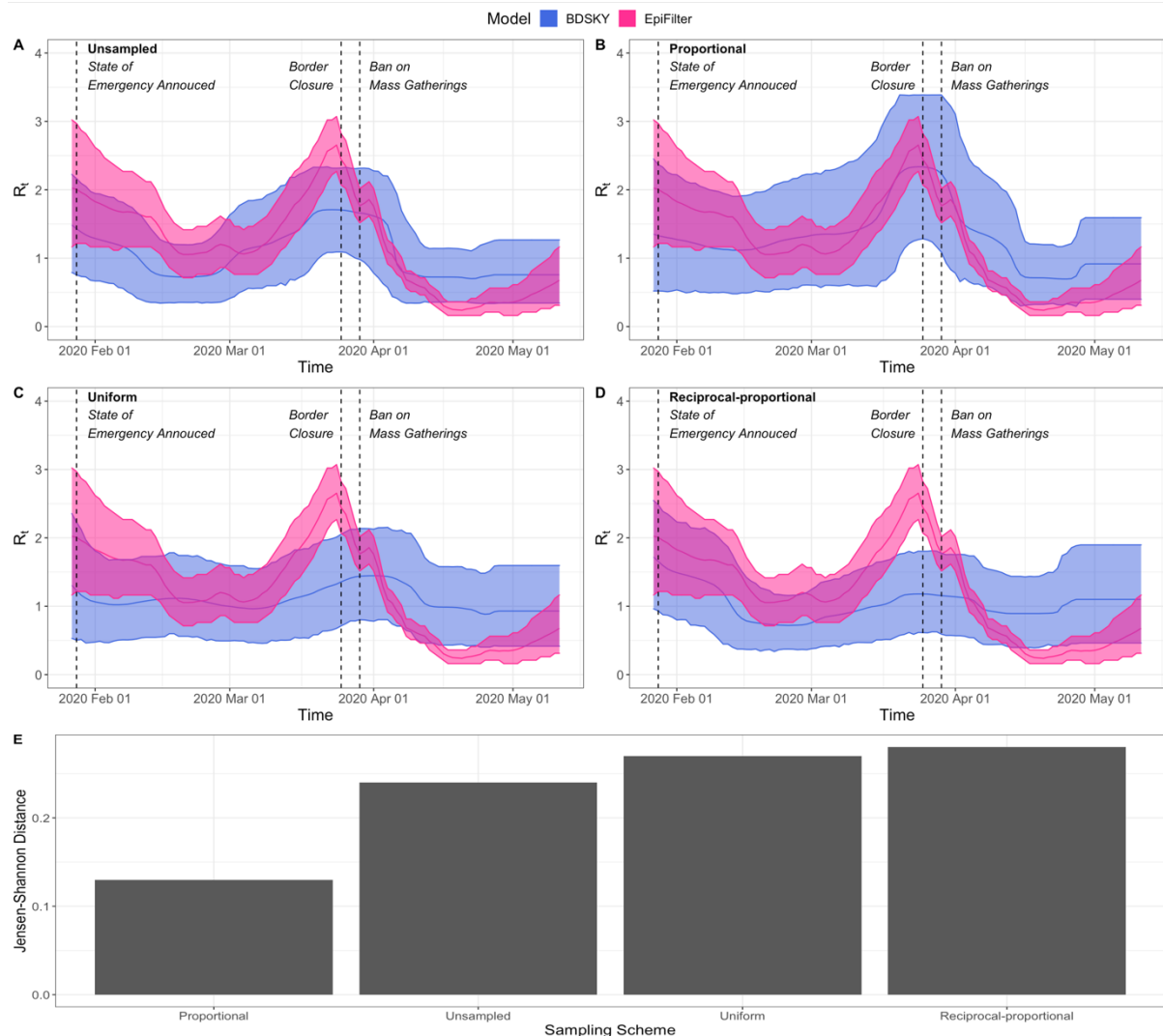
**Figure 3:**  $R_0$  estimated from BDSKY and TMRCA for Hong Kong and Brazil. Figure 1A and B represent Hong Kong and Figure 1C and D represent the Amazonas.

## Time-varying Reproduction number and Growth rate

We examine the  $R_t$  and  $r_t$  estimated for local SARS-CoV-2 epidemics in Hong Kong and Amazonas, Brazil. Our main results showing these two parameters and JSD are in figures 4-8.

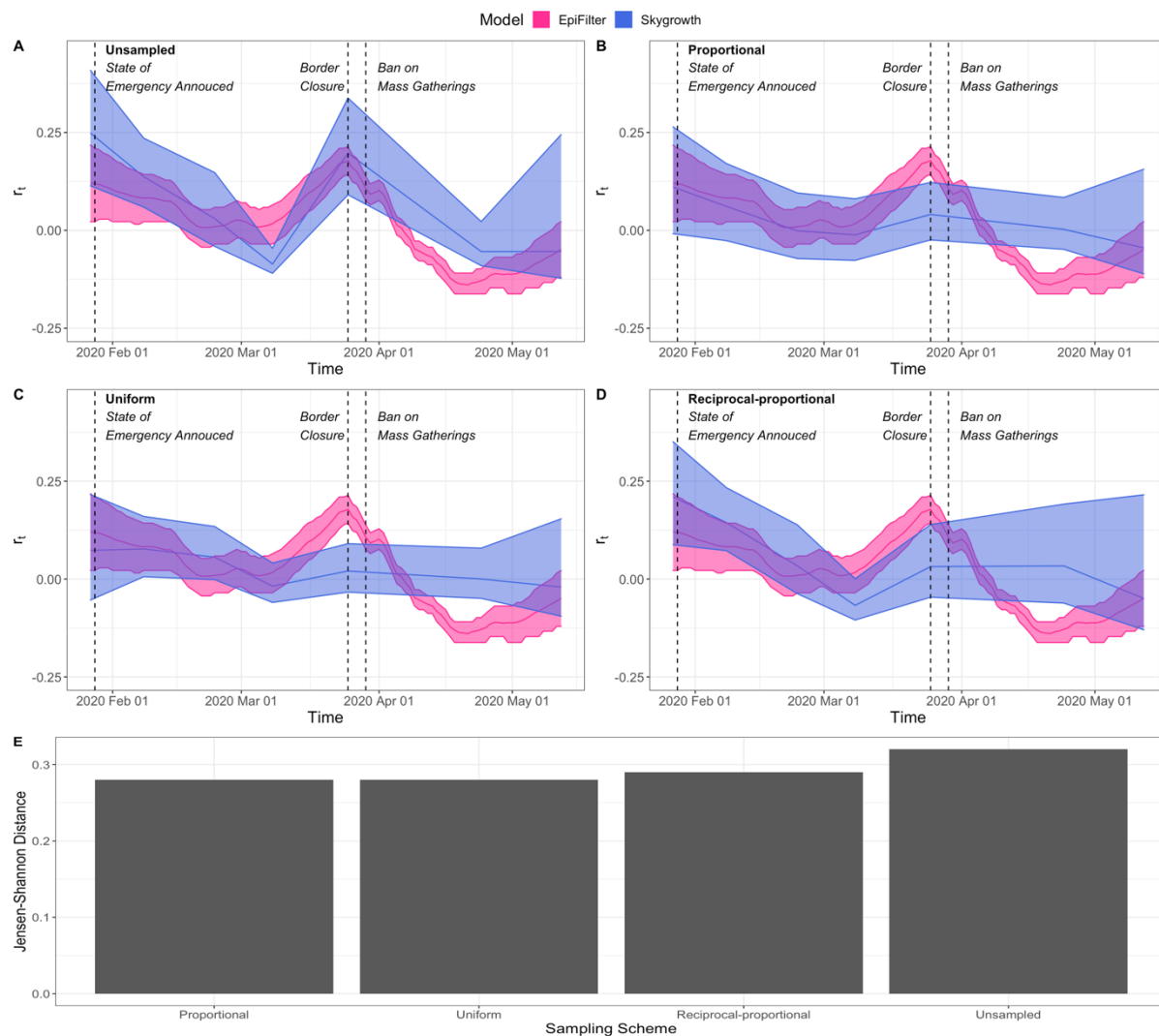
### *Hong Kong*

The BDSKY model was used alongside the EpiFilter model to estimate the  $R_t$  for each dataset subsampled according to the different sampling strategies (Figure 4). Based on the proportional sampling scheme, which had the lowest JSD (Figure 4E), we initially infer a super-critical  $R_t$  value, with a mean around an  $R_t$  value of 2, that appears to fall swiftly in response to the state of emergency and the rapid implementation of NPIs. A steady transmission rate subsequently persisted throughout the following weeks around the critical threshold ( $R_t = 1$ ). This period is succeeded by a sharp increase in  $R_t$ , peaking at a mean  $R_t$  value of 2.6. This is likely due to imported cases from North America and Europe (Cowling *et al.*, 2020). This led to a ban on international travel resulting in a sharp decline in  $R_t$  (Figure 2). However, this decline lasted around a week with the mean  $R_t$  briefly increasing until more stringent NPIs such as the banning of major gatherings were implemented. Following this, the  $R_t$  continued its sharp decline falling below the critical threshold, with transmission becoming sub-critical (Figure 4).



**Figure 4:**  $R_t$  estimated from both the BDSKY and *EpiFilter* models and Jensen Shannon Distance for Hong Kong. The bold writing represents the sampling scheme used in figure A-D. The light-shaded area represents the 95% HPDI with the darker-shaded area presenting where the BDSKY and *EpiFilter* models overlap. The solid line represents the mean  $R_t$  with *EpiFilter* being represented by a red line and BDSKY a blue line. The dashed lines represent policy change-times. The Jensen Shannon Distance is ordered from best to worse.

These results were mirrored in the estimation of  $r_t$  (Figure 5) for which the uniform and proportional sampling schemes showed the least divergence (Figure 5E). There was an initial decline in the  $r_t$ , which steadied at a value of  $\sim 0$ , indicating that epidemic stabilisation has occurred. This stable period is followed by an increase in  $r_t$  peaking at around a 5% increase in case incidence per day (Figure 5). In response to NPIs, the  $r_t$  starts to decrease, falling below 0, indicating a receding epidemic. The rate of this decline peaks at around a 7.5% reduction in case incidence per day (Figure 5).

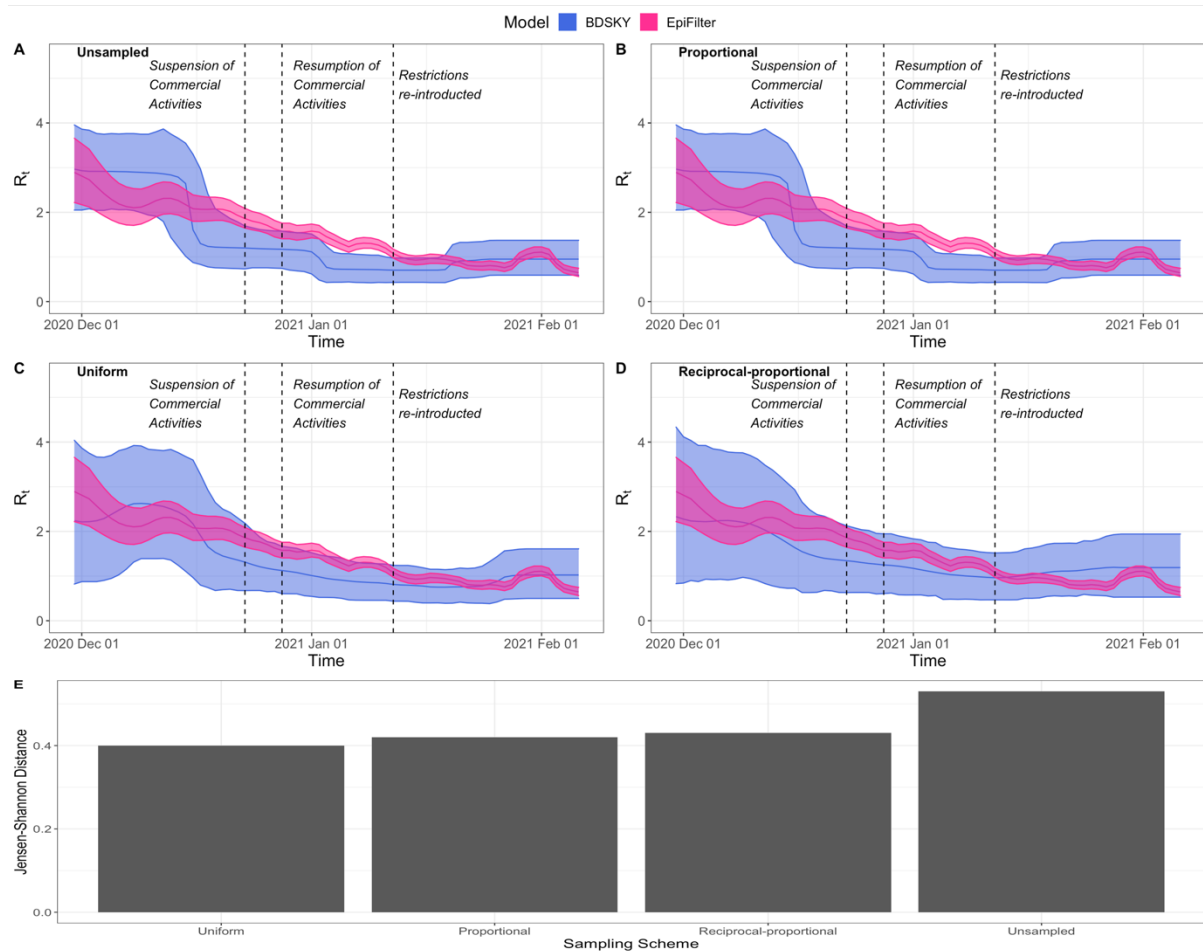


**Figure 5:**  $r_t$  estimated from both the *Skygrowth* and *EpiFilter* models and Jensen Shannon Distance for Hong Kong. The bold writing represents the sampling scheme used. The light-shaded area represents the 95% HPDI with the darker-shaded area presenting where the *Skygrowth* and *EpiFilter* models overlap. The solid line represents the mean  $r_t$  with *EpiFilter* being represented by a red line and *Skygrowth* a blue line. The dashed lines represent policy change-times. The Jensen Shannon Distance is ordered from best to worse.

### Brazil

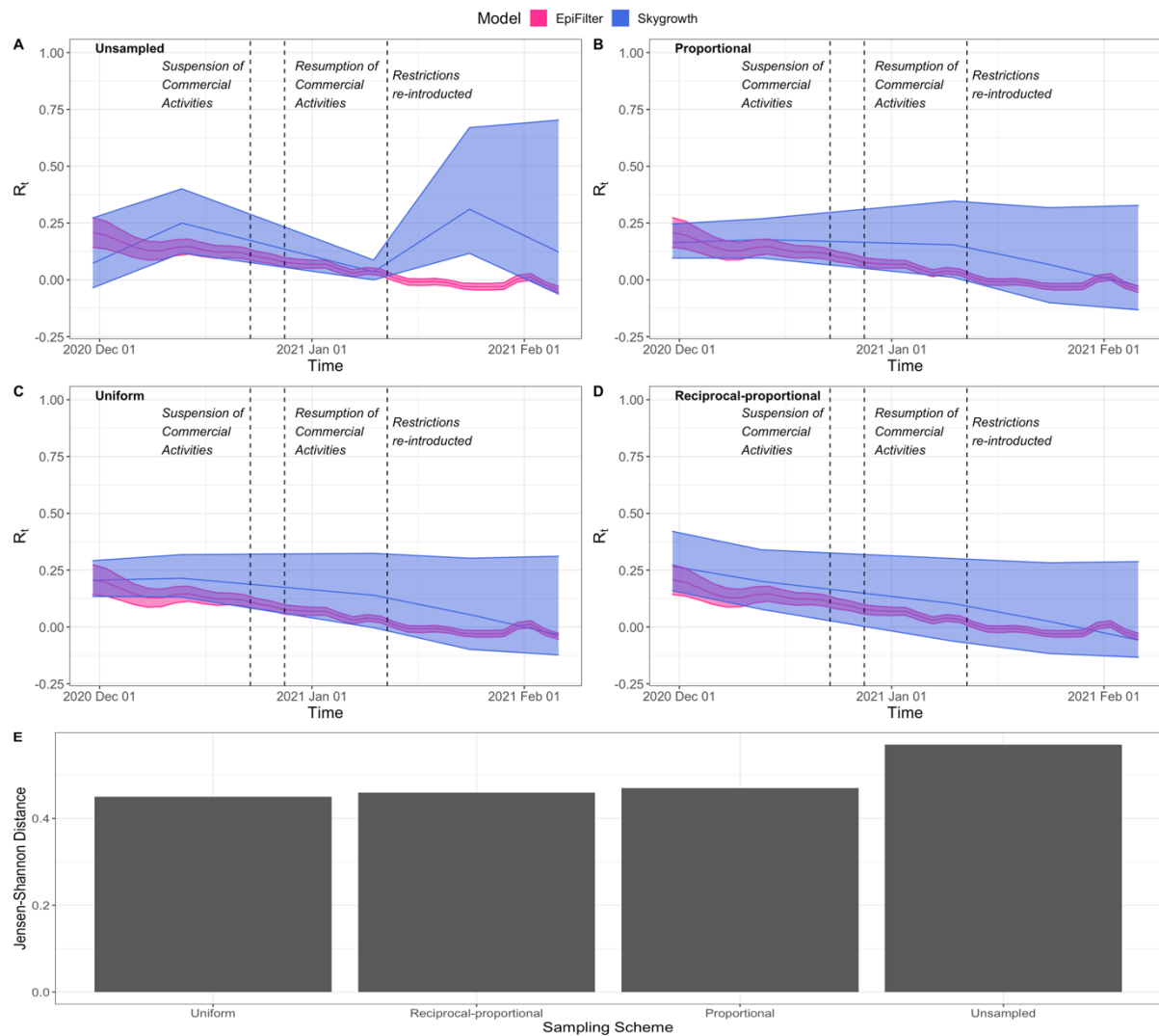
Based on the uniform sampling scheme, which had the lowest JSD (Figure 6E), we initially infer a super-critical  $R_t$  ( $R_t > 1$ ) value with a mean value of  $R_t = 3$  (Figure 6). From this point, the  $R_t$  declines, although it remains above the critical threshold ( $R_t = 1$ ) for much of the study period. Sub-critical ( $R_t < 1$ ) transmission was only reached after the re-imposition of NPIs. This implies that initial restrictions, such as the suspension of commercial activities, were ineffective in lowering the  $R_t$  below its critical threshold. Only after more stringent

restrictions were imposed did  $R_t$  become sub-critical. However, there is no evidence of a sharp decrease in  $R_t$  once restrictions were re-imposed, indicating they may have not had a significant impact on  $R_t$ .



**Figure 6:**  $R_t$  estimated from both the BDSKY and *EpiFilter* models and Jensen Shannon Distance for Amazonas, Brazil. The bold writing represents the sampling scheme used. The light-shaded area represents the 95% HPDI with the darker-shaded area presenting where the BDSKY and *EpiFilter* models overlap. The solid line represents the mean  $R_t$  with *EpiFilter* being represented by a blue line and BDSKY a red line. The dashed lines represent policy change-times. The Jensen Shannon Distance is ordered from best to worse.

Based on the uniform sampling which had the lowest JSD (Figure 7E) we infer a steady decline in  $r_t$  which matches the pattern seen with the  $R_t$  value (Figure 7). The initial  $r_t$  implied a 23% mean increase in case incidence per day. Subsequently, the  $r_t$  falls over the study period.  $r_t$  falls below 0 after the re-imposition of NPIs with a 3% reduction in mean case incidence per day by the end of the study period. There is no evidence of any noticeable declines in  $r_t$  when interventions were introduced indicating that they may have had a minimal impact on the growth rate of  $P.1/\gamma$ .



**Figure 7:**  $r_t$  estimated from both the *Skygrowth* and *EpiFilter* models and Jensen Shannon Distance for Amazonas, Brazil. The bold writing represents the sampling scheme used. The light-shaded area represents the 95% HPDI with the darker-shaded area presenting where the *EpiFilter* and *Skygrowth* models overlap. The solid line represents the mean  $r_t$  with *EpiFilter* being represented by a red line and *Skygrowth* a blue line. The dashed lines represent policy change-times. The Jensen Shannon Distance is ordered from best to worse.



## Discussion

In this study, phylodynamic methods have been applied to available SARS-CoV-2 sequences from Hong Kong and the Amazonas region of Brazil to infer their relevant epidemiological parameters and to compare the impact that various sampling strategies have on the phylodynamic reconstruction of these parameters.

We estimated the basic reproductive number of SARS-CoV-2 in Hong Kong to be 2.17 (95% CI = 1.43-2.83). This supports previous estimates of the initial  $R_0$  in Hong Kong (Cowling *et al.*, 2020; Zhao *et al.*, 2020) which estimates  $R_0$  to be 2.23 (95% CI = 1.47-3.42). For the Amazonas region in Brazil, we estimated the  $R_0$  to be 3.67 (95% CI = 2.83 – 4.48). Whilst the population of Amazonas State may not be fully susceptible to P.1/gamma (Faria *et al.*, 2021), this shouldn't affect the comparison between sampling schemes. Comparisons of different sampling schemes have revealed the  $R_0$  is robust to changes in sampling schemes (Figure 3A and C).

For the Hong Kong dataset, the proportional sampling scheme was superior to all other sampling schemes in estimating  $R_t$ . It successfully predicted the initial super-critical  $R_t$ , its decline in response to rapid NPIs, and subsequent increase and decline during the second wave of infections (Figure 4B). This was in comparison to the reciprocal-proportional that provided the worst JSD (Figure 4D) and in which the  $R_t$  remained relatively constant throughout the period. In addition, the proportional sampling scheme, alongside the uniform sampling scheme, best estimated  $r_t$  (Figure 5B and C). In contrast, for the Amazonas dataset, the uniform sampling scheme best estimated the  $R_t$  and was joint best for  $r_t$  (Figure 6C and Figure 7C). It captured both its initial super-critical  $R_t$  and high  $r_t$  alongside their subsequent decline. Our estimations for  $R_t$  are consistent with previous estimates of P.1 in Amazonas state (Faria *et al.*, 2021). This contrasted with the unsampled data in which the  $r_t$  increased at the end of the period (Figure 7A). This highlights that unlike  $R_0$ , both  $R_t$  and  $r_t$  are sensitive to changes in sampling and that even related epidemiological parameters like  $R_t$  and  $r_t$  may require different sampling strategies to optimise inferences.

Molecular clock dating of the Hong Kong and Amazonas dataset has revealed that the date of origin is robust to changes in sampling schemes. For Hong Kong, SARS-CoV-2 likely emerged in mid-December 2019 around 5 weeks before the first reported case on the 22<sup>nd</sup> of January 2020 (Cowling *et al.*, 2020). The Amazonas dataset revealed that the date of the

common ancestor of the P.1 lineage emerged around late October 2020, around 5 weeks before the first reported case on the 6<sup>th</sup> of December (Faria *et al.*, 2021), with all BCI's overlapping for each sampling strategy. Like the molecular clock dating, we found that the molecular clock rate was robust to changes in sampling strategies in both datasets with all sampling strategies having overlapped BCI's (Supplementary Table 2 and Supplementary Figure 4). For the Hong Kong dataset, its clock rate is comparable to early estimations of the mean substitution rate per site per year of SARS-CoV-2 (Duchene *et al.*, 2020). However, the clock rate estimated for the Brazilian dataset is lower than initial  $8.00 \times 10^{-4}$  s/s/y which is used in investigating SARS-CoV-2 (Andersen *et al.*, 2020) and that has been used in previous analyses of P.1 (Naveca *et al.*, 2021). This initial estimation of evolutionary rate was estimated from genomic data taken over a short time span at the beginning of the pandemic introducing a time dependency bias (Ghafari *et al.*, 2022). By using a more appropriate clock rate it can improve tree height and rooting resulting in more robust parameter estimations (Boskova, Stadler and Magnus, 2018).

Treating sampling times as uninformative has been shown to be inferior to including them as dependent on effective population size and other parameters by several previous studies (Hall, Woolhouse and Rambaut, 2016; Karcher *et al.*, 2016; Liu *et al.*, 2020; Parag, du Plessis and Pybus, 2020). Whilst these studies did not consider the estimation of epidemiological parameters, they highlight the potential of systematic biases being introduced into the phylodynamic reconstruction by not using a sampling scheme or by assuming an incorrect model for how sampling schemes introduce information. This was supported by our results as phylodynamic inferences with no sampling strategy applied had the poorest performance for both Hong Kong and the Amazonas region. This implies that sampling has a significant impact on phylodynamic reconstruction, and that exploration of sampling strategies is needed to obtain the most robust parameter estimates.

While our results provide a rigorous underpinning and insight into the dynamics of SARS-CoV-2 and the impact of sampling strategies in the Amazonas region and Hong Kong, there are limitations. The *Skygrowth* and BDSKY models do not explicitly consider imports into their respective regions. This is particularly relevant for Hong Kong as most initial sequences from the region were sequenced from importation events (Adam *et al.*, 2020) which can introduce error into parameter estimation. However, as the epidemic expanded, more

infections were attributable to autochthonous transmission (Adam *et al.*, 2020), and the risk of error introduced by importation events decreased. Moreover, while sampling strategies can account for temporal variations in genomic sampling fractions there is currently no way to account for non-random sampling approaches in either the BDSKY or *Skygrowth* models (Vasylyeva *et al.*, 2020). It is unclear how network-based sampling may affect parameter estimates obtained through these models (Volz, Koelle and Bedford, 2013) presenting a key challenge in molecular and genetic epidemiology. Spatial heterogeneities were also not explored within this work. This represents the next key step in understanding the impact of sampling as spatial sampling schemes would allow the reconstruction of the dispersal dynamics and estimation of epidemic overdispersion ( $k$ ), a key epidemiological parameter.

This work has highlighted the impact and importance that applying temporal sampling strategies can have on phylodynamic reconstruction. Whilst more genomic datasets from a variety of countries and regions with different sampling intensities and proportions are needed to create a more generalisable sampling framework and to dissect any potential cofounders, it has been shown that genomic datasets with no sampling strategy applied can introduce significant uncertainty and biases in the estimation of epidemiological parameters. This finding identifies the need for more targeted attempts at performing genomic surveillance and epidemic analyses particularly in resource-poor settings which have a limited genomic capability.

**Role of the Funding Sources:** N.R.F. acknowledges support from Wellcome Trust and Royal Society Sir Henry Dale Fellowship (204311/Z/16/Z), Bill and Melinda Gates Foundation (INV-034540) and Medical Research Council-Sao Paulo Research Foundation (FAPESP) CADDE partnership award (MR/S0195/1 and FAPESP 18/14389-0) (<https://caddecentre.org>). K.V.P. acknowledges support from grant reference MR/R015600/1, jointly funded by the UK Medical Research Council (MRC) and the UK Department for International Development (DFID) and from the NIHR Health Protection Research Unit in Behavioural Science and Evaluation at University of Bristol.

**CRedit authorship contribution statement:** R.P.D.I, K.V.P and N.R.F conceived and designed the study, R.P.D.I wrote and performed the analyses. R.P.D.I wrote the manuscript which was edited and supervised by K.V.P and N.R.F. All authors have contributed to and approved the manuscript for submission.

## References

- Adam, D.C. *et al.* (2020) “Clustering and superspreading potential of SARS-CoV-2 infections in Hong Kong,” *Nature Medicine*, 26(11), pp. 1714–1719. doi:10.1038/s41591-020-1092-0.
- Andersen, K.G. *et al.* (2020) “The proximal origin of SARS-CoV-2,” *Nature Medicine* [Preprint]. doi:10.1038/s41591-020-0820-9.
- Anderson *et al.* (2020) “The Royal Society SET-C Reports. Reproduction number (R) and growth rate (r) of the COVID-19 epidemic in the UK: methods of estimation, data sources, causes of heterogeneity, and use as a guide in policy formulation [report unpublished],” *The Royal Society*, (August), pp. 1–86.
- Anisimova, M. *et al.* (2011) “Survey of branch support methods demonstrates accuracy, power, and robustness of fast likelihood-based approximation schemes,” *Systematic biology*. 2011/05/03, 60(5), pp. 685–699. doi:10.1093/sysbio/syr041.
- Ayres, D.L. *et al.* (2019) “BEAGLE 3: Improved Performance, Scaling, and Usability for a High-Performance Computing Library for Statistical Phylogenetics,” *Systematic Biology*, 68(6), pp. 1052–1061. doi:10.1093/sysbio/syz020.
- Boskova, V., Stadler, T. and Magnus, C. (2018) “The influence of phylodynamic model specifications on parameter estimates of the Zika virus epidemic,” *Virus Evolution*, 4(1). doi:10.1093/ve/vex044.
- Bouckaert, R. *et al.* (2019) “BEAST 2.5: An advanced software platform for Bayesian evolutionary analysis,” *PLOS Computational Biology*, 15(4), p. e1006650. Available at: <https://doi.org/10.1371/journal.pcbi.1006650>.
- Byrne, A.W. *et al.* (2020) “Inferred duration of infectious period of SARS-CoV-2: Rapid scoping review and analysis of available evidence for asymptomatic and symptomatic COVID-19 cases,” *BMJ Open*, 10(8), pp. 1–16. doi:10.1136/bmjopen-2020-039856.
- Cori, A. *et al.* (2013) “A New Framework and Software to Estimate Time-Varying Reproduction Numbers During Epidemics,” *American Journal of Epidemiology*, 178(9), pp. 1505–1512. doi:10.1093/aje/kwt133.
- Cowling, B.J. *et al.* (2020) “Impact assessment of non-pharmaceutical interventions against coronavirus disease 2019 and influenza in Hong Kong: an observational study,” *The Lancet Public Health*, 5(5), pp. e279–e288. doi:10.1016/S2468-2667(20)30090-6.
- Dolan, P.T., Whitfield, Z.J. and Andino, R. (2018) “Mapping the Evolutionary Potential of RNA Viruses,” *Cell Host and Microbe*, 23(4), pp. 435–446. doi:10.1016/j.chom.2018.03.012.

- Drummond, A.J. *et al.* (2003) “Measurably evolving populations,” *Trends in Ecology & Evolution*, 18(9), pp. 481–488. doi:[https://doi.org/10.1016/S0169-5347\(03\)00216-7](https://doi.org/10.1016/S0169-5347(03)00216-7).
- Drummond, A.J. *et al.* (2005) “Bayesian Coalescent Inference of Past Population Dynamics from Molecular Sequences,” *Molecular Biology and Evolution*, 22(5), pp. 1185–1192. doi:[10.1093/molbev/msi103](https://doi.org/10.1093/molbev/msi103).
- Duchene, S. *et al.* (2020) “Temporal signal and the phylodynamic threshold of SARS-CoV-2,” *Virus Evolution*, 6(2). doi:[10.1093/ve/veaa061](https://doi.org/10.1093/ve/veaa061).
- Dudas, G. *et al.* (2017) “Virus genomes reveal factors that spread and sustained the Ebola epidemic,” *Nature* [Preprint]. doi:[10.1038/nature22040](https://doi.org/10.1038/nature22040).
- European Centre for Disease Prevention and Control (2020) *Guidelines for the implementation of non-pharmaceutical interventions against COVID-19 Key messages General considerations on NPI to control COVID-19*.
- Faria, N.R. *et al.* (2017) “Establishment and cryptic transmission of Zika virus in Brazil and the Americas,” *Nature*, 546(7658), pp. 406–410. doi:[10.1038/nature22401](https://doi.org/10.1038/nature22401).
- Faria, N.R. *et al.* (2021) “Genomics and epidemiology of the P.1 SARS-CoV-2 lineage in Manaus, Brazil,” *Science*, 372(6544), pp. 815 LP – 821. doi:[10.1126/science.abh2644](https://doi.org/10.1126/science.abh2644).
- Flaxman, S. *et al.* (2020) “Estimating the effects of non-pharmaceutical interventions on COVID-19 in Europe,” *Nature*, 584(7820), pp. 257–261. doi:[10.1038/s41586-020-2405-7](https://doi.org/10.1038/s41586-020-2405-7).
- Fraser, C. (2007) “Estimating Individual and Household Reproduction Numbers in an Emerging Epidemic,” *PLoS ONE*. Edited by A. Galvani, 2(8), p. e758. doi:[10.1371/journal.pone.0000758](https://doi.org/10.1371/journal.pone.0000758).
- Frost, S.D.W. *et al.* (2015) “Eight challenges in phylodynamic inference,” *Epidemics*. 2014/09/16, 10, pp. 88–92. doi:[10.1016/j.epidem.2014.09.001](https://doi.org/10.1016/j.epidem.2014.09.001).
- Ghafari, M. *et al.* (2022) “Purifying selection determines the short-term time dependency of evolutionary rates in SARS-CoV-2 and pH1N1 influenza,” *Molecular Biology and Evolution*, p. msac009. doi:[10.1093/molbev/msac009](https://doi.org/10.1093/molbev/msac009).
- Gill, M.S. *et al.* (2013) “Improving Bayesian population dynamics inference: a coalescent-based model for multiple loci,” *Molecular biology and evolution*. 2012/11/22, 30(3), pp. 713–724. doi:[10.1093/molbev/mss265](https://doi.org/10.1093/molbev/mss265).
- Gorbalenya, A.E. *et al.* (2020) “The species Severe acute respiratory syndrome-related coronavirus: classifying 2019-nCoV and naming it SARS-CoV-2,” *Nature Microbiology*, 5(4), pp. 536–544. doi:[10.1038/s41564-020-0695-z](https://doi.org/10.1038/s41564-020-0695-z).

- Gostic, K.M. *et al.* (2020) “Practical considerations for measuring the effective reproductive number,  $R_t$ ,” *PLOS Computational Biology*, 16(12), p. e1008409. Available at: <https://doi.org/10.1371/journal.pcbi.1008409>.
- Grubaugh, N.D. *et al.* (2017) “Genomic epidemiology reveals multiple introductions of Zika virus into the United States,” *Nature*, 546(7658), pp. 401–405. doi:10.1038/nature22400.
- Hadfield, J. *et al.* (2018) “Nextstrain: real-time tracking of pathogen evolution,” *Bioinformatics*, 34(23), pp. 4121–4123. doi:10.1093/bioinformatics/bty407.
- Hall, M.D., Woolhouse, M.E.J. and Rambaut, A. (2016) “The effects of sampling strategy on the quality of reconstruction of viral population dynamics using Bayesian skyline family coalescent methods: A simulation study,” *Virus evolution*, 2(1), pp. vew003–vew003. doi:10.1093/ve/vew003.
- Harvey, W.T. *et al.* (2021) “SARS-CoV-2 variants, spike mutations and immune escape,” *Nature Reviews Microbiology*, 19(7), pp. 409–424. doi:10.1038/s41579-021-00573-0.
- Hasegawa, M., Kishino, H. and Yano, T. (1985) “Dating of the human-ape splitting by a molecular clock of mitochondrial DNA,” *Journal of Molecular Evolution*, 22(2), pp. 160–174. doi:10.1007/BF02101694.
- Hidano, A. and Gates, M.C. (2019) “Assessing biases in phylodynamic inferences in the presence of super-spreaders,” *Veterinary Research*, 50(1), p. 74. doi:10.1186/s13567-019-0692-5.
- Hill, V. and Baele, G. (2019) “Bayesian Estimation of Past Population Dynamics in BEAST 1.10 Using the Skygrid Coalescent Model,” *Molecular Biology and Evolution*, 36(11), pp. 2620–2628. doi:10.1093/molbev/msz172.
- IBGE (2020) *Population Projections*. Available at: <https://www.ibge.gov.br/en/statistics/social/population.html> (Accessed: July 25, 2021).
- Jombart, T. *et al.* (2014) “Bayesian Reconstruction of Disease Outbreaks by Combining Epidemiologic and Genomic Data,” *PLOS Computational Biology*, 10(1), pp. e1003457-. Available at: <https://doi.org/10.1371/journal.pcbi.1003457>.
- Kalyaanamoorthy, S. *et al.* (2017) “ModelFinder: fast model selection for accurate phylogenetic estimates,” *Nature Methods*, 14(6), pp. 587–589. doi:10.1038/nmeth.4285.
- Karcher, M.D. *et al.* (2016) “Quantifying and Mitigating the Effect of Preferential Sampling on Phylodynamic Inference,” *PLoS computational biology*, 12(3), pp. e1004789–e1004789. doi:10.1371/journal.pcbi.1004789.



- Katoh, K. *et al.* (2002) “MAFFT: a novel method for rapid multiple sequence alignment based on fast Fourier transform,” *Nucleic Acids Research*, 30(14), pp. 3059–3066. doi:10.1093/nar/gkf436.
- Kullback, S. and Leibler, R.A. (1951) “On Information and Sufficiency,” *The Annals of Mathematical Statistics*, 22(1), pp. 79–86. doi:10.1214/aoms/1177729694.
- Lesley, S. *et al.* (2021) “Track Omicron’s spread with molecular data,” *Science*, 374(6574), pp. 1454–1455. doi:10.1126/science.abn4543.
- Lin, J. (1991) “Divergence measures based on the Shannon entropy,” *IEEE Transactions on Information Theory*, 37(1), pp. 145–151. doi:10.1109/18.61115.
- Liu, Q. *et al.* (2020) “Population Genetics of SARS-CoV-2: Disentangling Effects of Sampling Bias and Infection Clusters,” *Genomics, Proteomics & Bioinformatics*, 18(6), pp. 640–647. doi:<https://doi.org/10.1016/j.gpb.2020.06.001>.
- McAloon, C. *et al.* (2020) “Incubation period of COVID-19: a rapid systematic review and meta-analysis of observational research,” *BMJ Open*, 10(8), p. e039652. doi:10.1136/bmjopen-2020-039652.
- Minh, B.Q. *et al.* (2020) “IQ-TREE 2: New Models and Efficient Methods for Phylogenetic Inference in the Genomic Era,” *Molecular Biology and Evolution*, 37(5), pp. 1530–1534. doi:10.1093/molbev/msaa015.
- Nadeau, S.A. *et al.* (2021) “The origin and early spread of SARS-CoV-2 in Europe,” *Proceedings of the National Academy of Sciences*, 118(9), p. e2012008118. doi:10.1073/pnas.2012008118.
- Nascimento, V.A. do *et al.* (2020) “Genomic and phylogenetic characterisation of an imported case of SARS-CoV-2 in Amazonas State, Brazil,” *Memórias do Instituto Oswaldo Cruz*, 115. doi:10.1590/0074-02760200310.
- Naveca, F.G. *et al.* (2021) “COVID-19 in Amazonas, Brazil, was driven by the persistence of endemic lineages and P.1 emergence,” *Nature Medicine*, 27(7), pp. 1230–1238. doi:10.1038/s41591-021-01378-7.
- Parag, K. v. (2021) “Improved estimation of time-varying reproduction numbers at low case incidence and between epidemic waves,” *PLOS Computational Biology*. Edited by S.L. Kosakovsky Pond, 17(9), p. e1009347. doi:10.1371/journal.pcbi.1009347.
- Parag, K. v, du Plessis, L. and Pybus, O.G. (2020) “Jointly Inferring the Dynamics of Population Size and Sampling Intensity from Molecular Sequences,” *Molecular Biology and Evolution*, 37(8), pp. 2414–2429. doi:10.1093/molbev/msaa016.

- Prete, C.A. *et al.* (2021) “Serial interval distribution of SARS-CoV-2 infection in Brazil,” *Journal of travel medicine*, 28(2), pp. 1–3. doi:10.1093/jtm/taaa115.
- Pullano, G. *et al.* (2021) “Underdetection of cases of COVID-19 in France threatens epidemic control,” *Nature*, 590(7844), pp. 134–139. doi:10.1038/s41586-020-03095-6.
- Rai, B., Shukla, A. and Dwivedi, L.K. (2021) “Estimates of serial interval for COVID-19: A systematic review and meta-analysis,” *Clinical epidemiology and global health*. 2020/08/26, 9, pp. 157–161. doi:10.1016/j.cegh.2020.08.007.
- Rambaut, A. *et al.* (2016) “Exploring the temporal structure of heterochronous sequences using TempEst (formerly Path-O-Gen),” *Virus evolution*, 2(1), pp. vew007–vew007. doi:10.1093/ve/vew007.
- Rambaut, A. *et al.* (2018) “Posterior Summarization in Bayesian Phylogenetics Using Tracer 1.7,” *Systematic biology*, 67(5), pp. 901–904. doi:10.1093/sysbio/syy032.
- Rambaut, A. *et al.* (2020) “A dynamic nomenclature proposal for SARS-CoV-2 lineages to assist genomic epidemiology,” *Nature Microbiology*, 5(11), pp. 1403–1407. doi:10.1038/s41564-020-0770-5.
- Romano, C.M. and Melo, F.L. (2021) “Genomic surveillance of SARS-CoV-2: A race against time,” *The Lancet Regional Health - Americas*, 0(0), p. 100029. doi:10.1016/j.lana.2021.100029.
- Sabino, E.C. *et al.* (2021) “Resurgence of COVID-19 in Manaus, Brazil, despite high seroprevalence,” *Lancet (London, England)*, 397(10273), pp. 452–455. doi:10.1016/S0140-6736(21)00183-5.
- Shu, Y. and McCauley, J. (2017) “GISAID: Global initiative on sharing all influenza data - from vision to reality,” *Euro surveillance : bulletin Europeen sur les maladies transmissibles = European communicable disease bulletin*, 22(13), p. 30494. doi:10.2807/1560-7917.ES.2017.22.13.30494.
- de Silva, E., Ferguson, N.M. and Fraser, C. (2012) “Inferring pandemic growth rates from sequence data,” *Journal of The Royal Society Interface*, 9(73), pp. 1797–1808. doi:10.1098/rsif.2011.0850.
- de Souza, W.M. *et al.* (2020) “Epidemiological and clinical characteristics of the COVID-19 epidemic in Brazil,” *Nature Human Behaviour*, 4(8), pp. 856–865. doi:10.1038/s41562-020-0928-4.
- Stack, J.C. *et al.* (2010) “Protocols for sampling viral sequences to study epidemic dynamics,” *Journal of the Royal Society, Interface*. 2010/02/10, 7(48), pp. 1119–1127. doi:10.1098/rsif.2009.0530.



- Stadler, T. *et al.* (2013) “Birth–death skyline plot reveals temporal changes of epidemic spread in HIV and hepatitis C virus (HCV),” *Proceedings of the National Academy of Sciences*, 110(1), pp. 228 LP – 233. doi:10.1073/pnas.1207965110.
- Suchard, M.A. *et al.* (2018) “Bayesian phylogenetic and phylodynamic data integration using BEAST 1.10,” *Virus evolution*, 4(1), pp. vey016–vey016. doi:10.1093/ve/vey016.
- Tamura, K. and Nei, M. (1993) “Estimation of the number of nucleotide substitutions in the control region of mitochondrial DNA in humans and chimpanzees.,” *Molecular Biology and Evolution*, 10(3), pp. 512–526. doi:10.1093/oxfordjournals.molbev.a040023.
- The World Bank (2021) *Population, total - Hong Kong SAR, China*. Available at: <https://data.worldbank.org/indicator/SP.POP.TOTL?locations=HK> (Accessed: July 25, 2021).
- Tsang, T.K. *et al.* (2020) “Effect of changing case definitions for COVID-19 on the epidemic curve and transmission parameters in mainland China: a modelling study.,” *The Lancet. Public health*, 5(5), pp. e289–e296. doi:10.1016/S2468-2667(20)30089-X.
- UK Health Security Agency (2022) *The R value and growth rate*, <https://www.gov.uk/guidance/the-r-value-and-growth-rate>.
- Vasylyeva, T.I. *et al.* (2020) “Phylodynamics helps to evaluate the impact of an HIV prevention intervention,” *Viruses*, 12(4), pp. 1–15. doi:10.3390/v12040469.
- Verity, R. *et al.* (2020) “Estimates of the severity of coronavirus disease 2019: a model-based analysis,” *The Lancet. Infectious diseases*, 20(6), pp. 669–677. doi:10.1016/S1473-3099(20)30243-7.
- Volz, E. *et al.* (2021) “Evaluating the Effects of SARS-CoV-2 Spike Mutation D614G on Transmissibility and Pathogenicity,” *Cell*, 184(1), pp. 64–75.e11. doi:<https://doi.org/10.1016/j.cell.2020.11.020>.
- Volz, E.M. and Didelot, X. (2018) “Modeling the Growth and Decline of Pathogen Effective Population Size Provides Insight into Epidemic Dynamics and Drivers of Antimicrobial Resistance,” *Systematic Biology*, 67(4), pp. 719–728. doi:10.1093/sysbio/syy007.
- Volz, E.M., Koelle, K. and Bedford, T. (2013) “Viral phylodynamics,” *PLoS computational biology*. 2013/03/21, 9(3), pp. e1002947–e1002947. doi:10.1371/journal.pcbi.1002947.
- Wallinga, J. and Teunis, P. (2004) “Different Epidemic Curves for Severe Acute Respiratory Syndrome Reveal Similar Impacts of Control Measures,” *American Journal of Epidemiology*, 160(6), pp. 509–516. doi:10.1093/aje/kwh255.

Wallinga and Lipsitch (2007) “How generation intervals shape the relationship between growth rates and reproductive numbers,” *Proceedings of the Royal Society B: Biological Sciences*, 274(1609), pp. 599–604. doi:10.1098/rspb.2006.3754.

World Health Organisation (2020) *Public Health Emergency of International Concern (PHEIC)*.

World Health Organisation (2021a) *Genomic sequencing of SARS-CoV-2 A guide to implementation for maximum impact on public health*.

World Health Organisation (2021b) *Guidance for surveillance of SARS-CoV-2 variants Interim guidance*.

World Health Organisation (2022) *Coronavirus disease (COVID-19) Weekly Epidemiological Update and Weekly Operational Update*, <https://www.who.int/emergencies/diseases/novel-coronavirus-2019/situation-reports>.

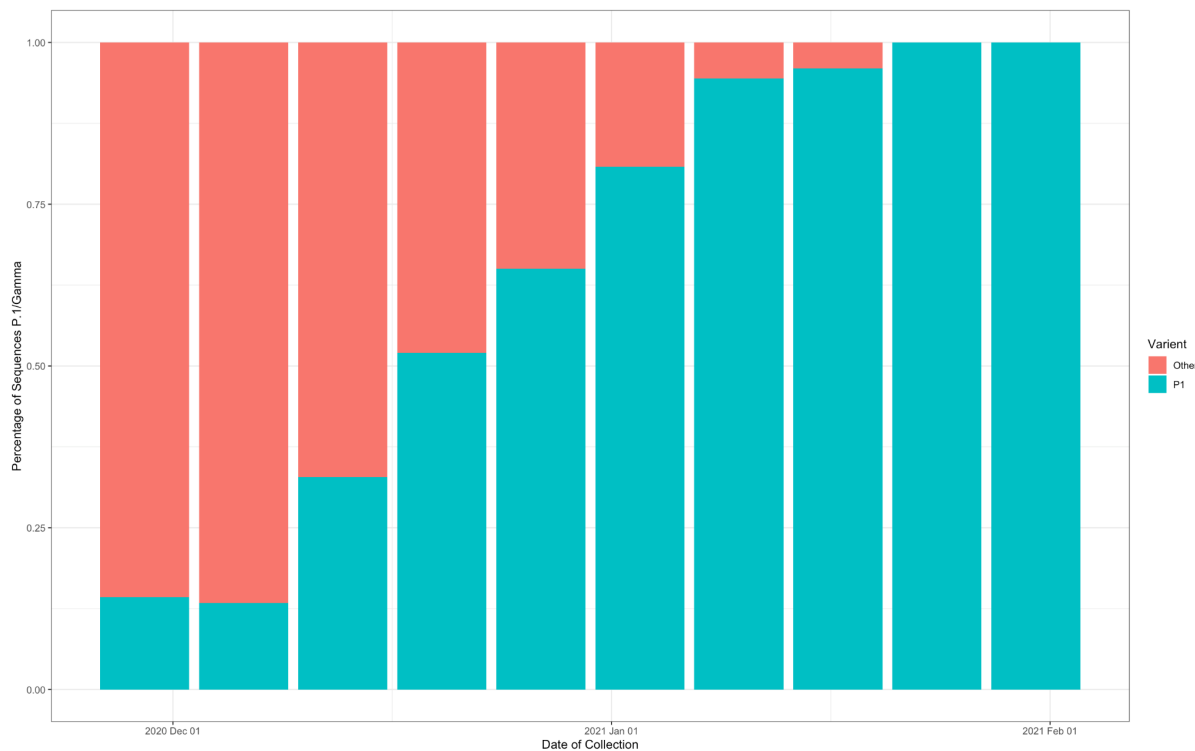
Zhao, S. *et al.* (2020) “Preliminary estimation of the basic reproduction number of novel coronavirus (2019-nCoV) in China, from 2019 to 2020: A data-driven analysis in the early phase of the outbreak,” *International journal of infectious diseases : IJID : official publication of the International Society for Infectious Diseases*. 2020/01/30, 92, pp. 214–217. doi:10.1016/j.ijid.2020.01.050.

Zhu, N. *et al.* (2020) “A Novel Coronavirus from Patients with Pneumonia in China, 2019,” *New England Journal of Medicine*, 382(8), pp. 727–733. doi:10.1056/NEJMoa2001017.

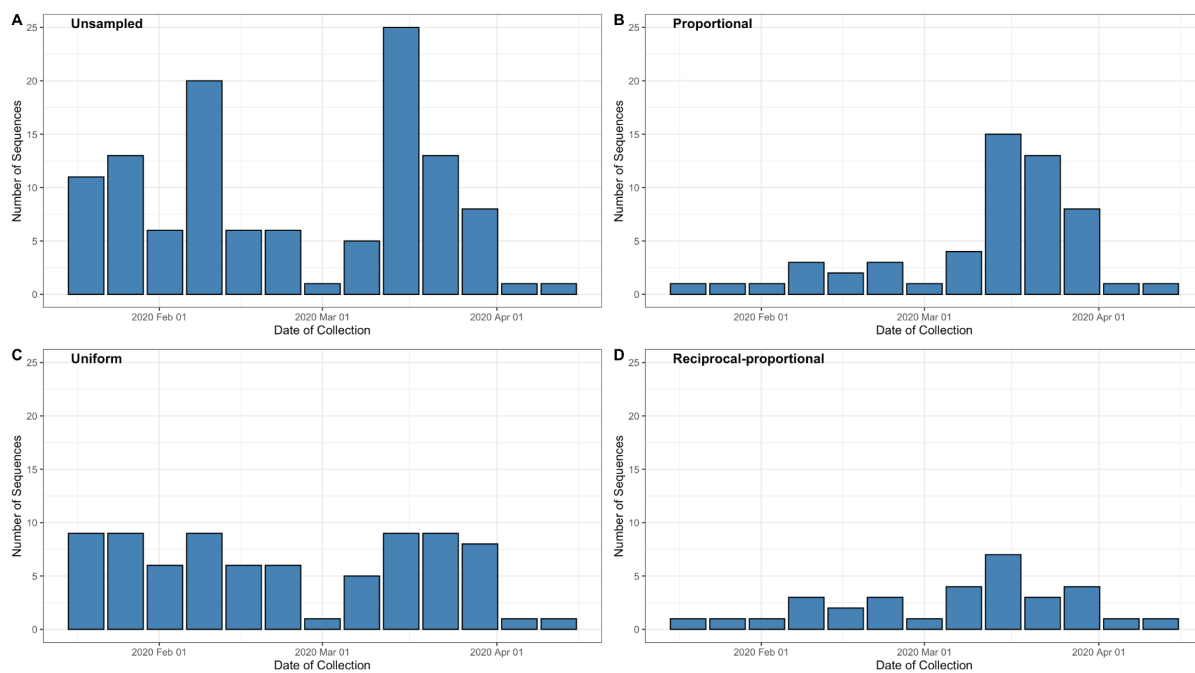
## Supplementary Figures and Tables

**Supplementary Table 1:** Key parameters and definitions for SARS-CoV-2

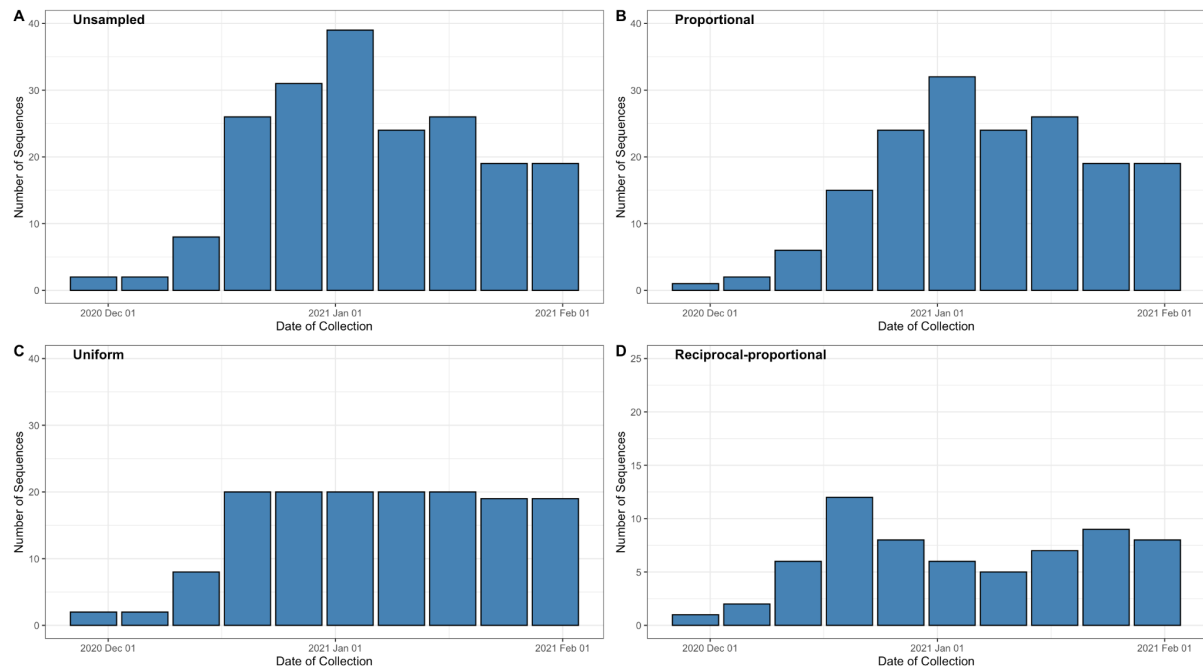
Parameter	Definition
Basic reproduction number ( $R_0$ )	Average number of individuals infected by a single infected person in a fully susceptible population
Time-varying or effective reproduction number ( $R_t$ )	Average number of secondary infections generated per effective primary case at a certain time point and in the presence of susceptible depletion or interventions
Growth rate ( $r_t$ )	Rate of change of the logarithm of the number of new cases per unit of time
Incubation period	Time between infection and symptom onset
Infectious period	Period in which an infectious host can transmit infectious agents to a susceptible individual
Generation interval	Time between infection events in an infector–infectee pair
Date of origin	Date in which viral variant is thought to have emerged
Serial Interval	Time between symptom onsets in an infector–infectee pair



**Supplementary Figure 1:** The proportion of P.1 sequences compared to non-P.1 sequences found on GISAID (Shu and McCauley, 2017).



**Supplementary Figure 2:** Number of sequences for each week and sampling scheme for Hong Kong dataset.

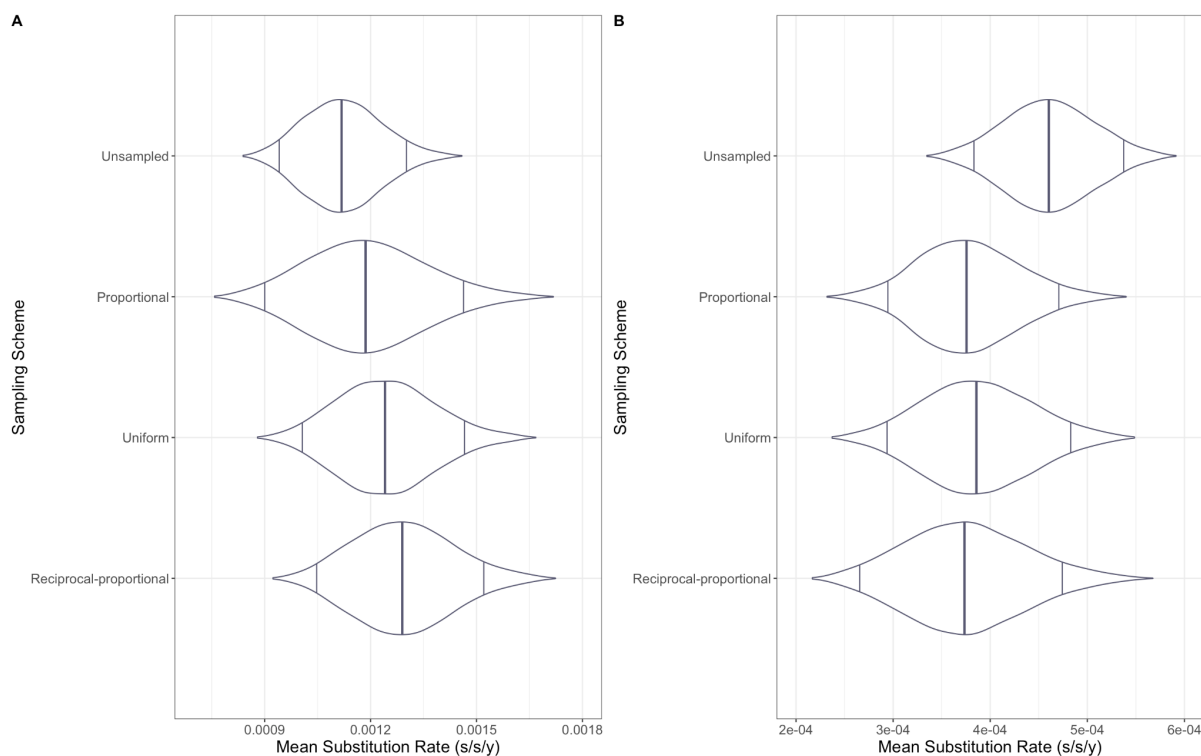


**Supplementary Figure 3:** Number of sequences for each week and sampling scheme for Amazonas dataset.

**Supplementary Table 2:** TMRCA and mean substitution rate both with 95% BCI for each sampling strategy for Hong Kong and Amazonas datasets alongside the Jensen-Shannon distance. Full posterior distribution of the TMRCA and substitution rates obtained under the different sampling strategies can be found in Figure 3B and D and Supplementary Figure 4.

Sampling Strategy	Dataset	TMRCA (95% BCI)	Mean Substitution Rate (95% BCI, subs/site/year, s/s/y)
Unsampled	Hong Kong	2 <sup>nd</sup> December 2019 (10 <sup>th</sup> November 2019 – 24 <sup>th</sup> December 2019)	$1.12 \times 10^{-3}$ ( $9.16 \times 10^{-4}$ – $1.35 \times 10^{-3}$ )
	Brazil	30 <sup>th</sup> October 2020 (8 <sup>th</sup> October 2020 – 13 <sup>th</sup> December 2020)	$4.58 \times 10^{-4}$ ( $3.69 \times 10^{-4}$ – $5.56 \times 10^{-4}$ )

Proportional	Hong Kong	24 <sup>th</sup> December 2019 (21 <sup>st</sup> November 2019 – 11 <sup>th</sup> January 2020)	$1.39 \times 10^{-3}$ ( $9.28 \times 10^{-4}$ – $2.48 \times 10^{-3}$ )
	Brazil	30 <sup>th</sup> October 2020 (25 <sup>th</sup> August 2020 – 29 <sup>th</sup> November 2020)	$4.60 \times 10^{-4}$ ( $3.70 \times 10^{-4}$ – $5.56 \times 10^{-4}$ )
Uniform	Hong Kong	13 <sup>th</sup> December 2019 (18 <sup>th</sup> November 2019 – 4 <sup>th</sup> January 2020)	$1.64 \times 10^{-3}$ ( $1.22 \times 10^{-3}$ – $2.09 \times 10^{-3}$ )
	Brazil	27 <sup>th</sup> October 2020 (5 <sup>th</sup> October 2020 – 25 <sup>th</sup> November 2020)	$4.60 \times 10^{-4}$ ( $3.70 \times 10^{-4}$ – $5.56 \times 10^{-4}$ )
Reciprocal-proportional	Hong Kong	6 <sup>th</sup> December 2019 (10 <sup>th</sup> November 2019 – 28 <sup>th</sup> December 2019)	$1.30 \times 10^{-3}$ ( $1.03 \times 10^{-3}$ – $1.59 \times 10^{-3}$ )
	Brazil	30 <sup>th</sup> October 2020 (27 <sup>th</sup> September 2020 – 25 <sup>th</sup> November 2020)	$4.00 \times 10^{-4}$ ( $2.56 \times 10^{-4}$ – $5.55 \times 10^{-4}$ )



**Supplementary Figure 4:** Mean substitution rate (s/s/y) for Hong Kong and Brazil. Figure 1A represents Hong Kong with Figure 1B representing the Amazonas.

**Supplementary Table 3:** Accession ID of each Hong Kong sequence for each sampling strategy used within this study

Unsampled	Proportional	Uniform	Reciprocal-proportional
EPI_ISL_412028	EPI_ISL_414517	EPI_ISL_412029	EPI_ISL_412028
EPI_ISL_412029	EPI_ISL_414519	EPI_ISL_414517	EPI_ISL_412029
EPI_ISL_412030	EPI_ISL_414527	EPI_ISL_414519	EPI_ISL_412030
EPI_ISL_414517	EPI_ISL_418815	EPI_ISL_414527	EPI_ISL_414517
EPI_ISL_414519	EPI_ISL_419224	EPI_ISL_414569	EPI_ISL_414519
EPI_ISL_414527	EPI_ISL_419229	EPI_ISL_414571	EPI_ISL_414527
EPI_ISL_414528	EPI_ISL_419232	EPI_ISL_416314	EPI_ISL_414528
EPI_ISL_414569	EPI_ISL_450404	EPI_ISL_417064	EPI_ISL_414569
EPI_ISL_414571	EPI_ISL_450405	EPI_ISL_417443	EPI_ISL_414571
EPI_ISL_416314	EPI_ISL_450410	EPI_ISL_419214	EPI_ISL_416314

EPI_ISL_ 417064	EPI_ISL_ 476801	EPI_ISL_ 419215	EPI_ISL_ 417064
EPI_ISL_ 417176	EPI_ISL_ 476802	EPI_ISL_ 419217	EPI_ISL_ 417176
EPI_ISL_ 417178	EPI_ISL_ 476803	EPI_ISL_ 419224	EPI_ISL_ 417178
EPI_ISL_ 417181	EPI_ISL_ 497769	EPI_ISL_ 419225	EPI_ISL_ 417181
EPI_ISL_ 417185	EPI_ISL_ 497773	EPI_ISL_ 419227	EPI_ISL_ 417185
EPI_ISL_ 417187	EPI_ISL_ 497775	EPI_ISL_ 419228	EPI_ISL_ 417187
EPI_ISL_ 417188	EPI_ISL_ 497784	EPI_ISL_ 419229	EPI_ISL_ 417188
EPI_ISL_ 417193	EPI_ISL_ 497786	EPI_ISL_ 419231	EPI_ISL_ 417193
EPI_ISL_ 417197	EPI_ISL_ 497791	EPI_ISL_ 419232	EPI_ISL_ 417197
EPI_ISL_ 417443	EPI_ISL_ 497796	EPI_ISL_ 419245	EPI_ISL_ 417443
EPI_ISL_ 418815	EPI_ISL_ 497799	EPI_ISL_ 419247	EPI_ISL_ 418815
EPI_ISL_ 419214	EPI_ISL_ 497806	EPI_ISL_ 419250	EPI_ISL_ 419214
EPI_ISL_ 419215	EPI_ISL_ 497808	EPI_ISL_ 419252	EPI_ISL_ 419215
EPI_ISL_ 419216	EPI_ISL_ 497810	EPI_ISL_ 434564	EPI_ISL_ 419216
EPI_ISL_ 419217	EPI_ISL_ 497811	EPI_ISL_ 434565	EPI_ISL_ 419217
EPI_ISL_ 419219	EPI_ISL_ 497818	EPI_ISL_ 434567	EPI_ISL_ 419219
EPI_ISL_ 419221	EPI_ISL_ 497819	EPI_ISL_ 434568	EPI_ISL_ 419221
EPI_ISL_ 419222	EPI_ISL_ 497821	EPI_ISL_ 434569	EPI_ISL_ 419222
EPI_ISL_ 419224	EPI_ISL_ 497823	EPI_ISL_ 434570	EPI_ISL_ 419224
EPI_ISL_ 419225	EPI_ISL_ 497824	EPI_ISL_ 434571	EPI_ISL_ 419225
EPI_ISL_ 419226	EPI_ISL_ 497840	EPI_ISL_ 450405	EPI_ISL_ 419226
EPI_ISL_ 419227	EPI_ISL_ 497845	EPI_ISL_ 450408	EPI_ISL_ 419227
EPI_ISL_ 419228	EPI_ISL_ 497846	EPI_ISL_ 450409	EPI_ISL_ 419228
EPI_ISL_ 419229	EPI_ISL_ 497847	EPI_ISL_ 450410	EPI_ISL_ 419229
EPI_ISL_ 419231	EPI_ISL_ 497850	EPI_ISL_ 450411	EPI_ISL_ 419231
EPI_ISL_ 419232	EPI_ISL_ 497856	EPI_ISL_ 476801	EPI_ISL_ 419232
EPI_ISL_ 419245	EPI_ISL_ 497865	EPI_ISL_ 476802	EPI_ISL_ 419245
EPI_ISL_ 419247	EPI_ISL_ 497870	EPI_ISL_ 476804	EPI_ISL_ 419247
EPI_ISL_ 419250	EPI_ISL_ 516798	EPI_ISL_ 497769	EPI_ISL_ 419250
EPI_ISL_ 419252	EPI_ISL_ 539820	EPI_ISL_ 497771	EPI_ISL_ 419252
EPI_ISL_ 434560	EPI_ISL_ 539850	EPI_ISL_ 497783	EPI_ISL_ 434563
EPI_ISL_ 434563	EPI_ISL_ 539851	EPI_ISL_ 497784	EPI_ISL_ 434564



EPI_ISL_ 434564	EPI_ISL_ 610167	EPI_ISL_ 497791	EPI_ISL_ 434565
EPI_ISL_ 434565	EPI_ISL_ 610168	EPI_ISL_ 497806	EPI_ISL_ 434566
EPI_ISL_ 434566	EPI_ISL_ 610169	EPI_ISL_ 497810	EPI_ISL_ 434567
EPI_ISL_ 434567	EPI_ISL_ 610170	EPI_ISL_ 497811	EPI_ISL_ 434568
EPI_ISL_ 434568	EPI_ISL_ 610171	EPI_ISL_ 497813	EPI_ISL_ 434569
EPI_ISL_ 434569	EPI_ISL_ 610172	EPI_ISL_ 497818	EPI_ISL_ 434570
EPI_ISL_ 434570	EPI_ISL_ 610173	EPI_ISL_ 497821	EPI_ISL_ 434571
EPI_ISL_ 434571	EPI_ISL_ 610174	EPI_ISL_ 497823	EPI_ISL_ 450405
EPI_ISL_ 450404	EPI_ISL_ 610175	EPI_ISL_ 497824	EPI_ISL_ 450408
EPI_ISL_ 450405	EPI_ISL_ 610177	EPI_ISL_ 497826	EPI_ISL_ 450409
EPI_ISL_ 450408		EPI_ISL_ 497827	EPI_ISL_ 450410
EPI_ISL_ 450409		EPI_ISL_ 497831	EPI_ISL_ 450411
EPI_ISL_ 450410		EPI_ISL_ 497832	EPI_ISL_ 450412
EPI_ISL_ 450411		EPI_ISL_ 497846	EPI_ISL_ 476802
EPI_ISL_ 450412		EPI_ISL_ 497847	EPI_ISL_ 476804
EPI_ISL_ 476801		EPI_ISL_ 497848	EPI_ISL_ 497769
EPI_ISL_ 476802		EPI_ISL_ 497856	EPI_ISL_ 497771
EPI_ISL_ 476803		EPI_ISL_ 497860	EPI_ISL_ 497773
EPI_ISL_ 476804		EPI_ISL_ 497865	EPI_ISL_ 497783
EPI_ISL_ 497769		EPI_ISL_ 539820	EPI_ISL_ 497784
EPI_ISL_ 497771		EPI_ISL_ 539850	EPI_ISL_ 497791
EPI_ISL_ 497773		EPI_ISL_ 539851	EPI_ISL_ 497797
EPI_ISL_ 497775		EPI_ISL_ 610165	EPI_ISL_ 497811
EPI_ISL_ 497783		EPI_ISL_ 610166	EPI_ISL_ 497812
EPI_ISL_ 497784		EPI_ISL_ 610167	EPI_ISL_ 497818
EPI_ISL_ 497786		EPI_ISL_ 610168	EPI_ISL_ 497819
EPI_ISL_ 497791		EPI_ISL_ 610169	EPI_ISL_ 497823
EPI_ISL_ 497796		EPI_ISL_ 610171	EPI_ISL_ 497824

EPI_ISL_ 497797		EPI_ISL_ 610173	EPI_ISL_ 497827
EPI_ISL_ 497798		EPI_ISL_ 610174	EPI_ISL_ 497831
EPI_ISL_ 497799		EPI_ISL_ 610175	EPI_ISL_ 497833
EPI_ISL_ 497806		EPI_ISL_ 610177	EPI_ISL_ 497848
EPI_ISL_ 497808			EPI_ISL_ 497850
EPI_ISL_ 497810			EPI_ISL_ 497856
EPI_ISL_ 497811			EPI_ISL_ 497860
EPI_ISL_ 497812			EPI_ISL_ 497864
EPI_ISL_ 497813			EPI_ISL_ 497865
EPI_ISL_ 497818			EPI_ISL_ 539850
EPI_ISL_ 497819			EPI_ISL_ 539851
EPI_ISL_ 497820			EPI_ISL_ 610165
EPI_ISL_ 497821			EPI_ISL_ 610166
EPI_ISL_ 497823			EPI_ISL_ 610172
EPI_ISL_ 497824			EPI_ISL_ 610177
EPI_ISL_ 497826			
EPI_ISL_ 497827			
EPI_ISL_ 497831			
EPI_ISL_ 497832			
EPI_ISL_ 497833			
EPI_ISL_ 497840			
EPI_ISL_ 497845			
EPI_ISL_ 497846			
EPI_ISL_ 497847			
EPI_ISL_ 497848			
EPI_ISL_ 497850			
EPI_ISL_ 497856			

EPI_ISL_ 497860			
EPI_ISL_ 497864			
EPI_ISL_ 497865			
EPI_ISL_ 497870			
EPI_ISL_ 516798			
EPI_ISL_ 539820			
EPI_ISL_ 539850			
EPI_ISL_ 539851			
EPI_ISL_ 610165			
EPI_ISL_ 610166			
EPI_ISL_ 610167			
EPI_ISL_ 610168			
EPI_ISL_ 610169			
EPI_ISL_ 610170			
EPI_ISL_ 610171			
EPI_ISL_ 610172			
EPI_ISL_ 610173			
EPI_ISL_ 610174			
EPI_ISL_ 610175			
EPI_ISL_ 610177			

**Supplementary Table 4:** Accession ID of each Amazonas State, Brazil sequence for each sampling strategy used within this study

<b>Unsampled</b>	<b>Proportional</b>	<b>Uniform</b>	<b>Reciprocal-proportional</b>
EPI_ISL_1034306	EPI_ISL_1034304	EPI_ISL_1034304	EPI_ISL_1034306
EPI_ISL_1060876	EPI_ISL_1034306	EPI_ISL_1034306	EPI_ISL_1060913
EPI_ISL_1060877	EPI_ISL_1060877	EPI_ISL_1060877	EPI_ISL_1060914
EPI_ISL_1060881	EPI_ISL_1060881	EPI_ISL_1060881	EPI_ISL_1068149
EPI_ISL_1060888	EPI_ISL_1060897	EPI_ISL_1060888	EPI_ISL_1068150
EPI_ISL_1060889	EPI_ISL_1060900	EPI_ISL_1060889	EPI_ISL_1068156
EPI_ISL_1060894	EPI_ISL_1060902	EPI_ISL_1060897	EPI_ISL_1068198
EPI_ISL_1060897	EPI_ISL_1060904	EPI_ISL_1060900	EPI_ISL_1068258
EPI_ISL_1060900	EPI_ISL_1060906	EPI_ISL_1060912	EPI_ISL_1068260
EPI_ISL_1060902	EPI_ISL_1060912	EPI_ISL_1060913	EPI_ISL_1068262
EPI_ISL_1060904	EPI_ISL_1060913	EPI_ISL_1060956	EPI_ISL_1068263
EPI_ISL_1060906	EPI_ISL_1060914	EPI_ISL_1061026	EPI_ISL_1068264
EPI_ISL_1060911	EPI_ISL_1060918	EPI_ISL_1068111	EPI_ISL_1068278
EPI_ISL_1060912	EPI_ISL_1060956	EPI_ISL_1068149	EPI_ISL_1068286
EPI_ISL_1060913	EPI_ISL_1061026	EPI_ISL_1068150	EPI_ISL_1068288
EPI_ISL_1060914	EPI_ISL_1068110	EPI_ISL_1068154	EPI_ISL_1166615
EPI_ISL_1060918	EPI_ISL_1068111	EPI_ISL_1068158	EPI_ISL_1213190
EPI_ISL_1060956	EPI_ISL_1068112	EPI_ISL_1068160	EPI_ISL_1261690
EPI_ISL_1061026	EPI_ISL_1068114	EPI_ISL_1068169	EPI_ISL_1261694
EPI_ISL_1068110	EPI_ISL_1068149	EPI_ISL_1068198	EPI_ISL_2777236
EPI_ISL_1068111	EPI_ISL_1068150	EPI_ISL_1068222	EPI_ISL_2777320
EPI_ISL_1068112	EPI_ISL_1068151	EPI_ISL_1068225	EPI_ISL_2777363
EPI_ISL_1068114	EPI_ISL_1068154	EPI_ISL_1068226	EPI_ISL_2777375
EPI_ISL_1068149	EPI_ISL_1068156	EPI_ISL_1068243	EPI_ISL_2777376
EPI_ISL_1068150	EPI_ISL_1068158	EPI_ISL_1068248	EPI_ISL_2777384
EPI_ISL_1068151	EPI_ISL_1068160	EPI_ISL_1068249	EPI_ISL_2777388
EPI_ISL_1068154	EPI_ISL_1068169	EPI_ISL_1068260	EPI_ISL_2777397

EPI_ISL_1068156	EPI_ISL_1068198	EPI_ISL_1068261	EPI_ISL_2777399
EPI_ISL_1068158	EPI_ISL_1068221	EPI_ISL_1068262	EPI_ISL_2777401
EPI_ISL_1068160	EPI_ISL_1068222	EPI_ISL_1068263	EPI_ISL_2777403
EPI_ISL_1068169	EPI_ISL_1068225	EPI_ISL_1068264	EPI_ISL_2777404
EPI_ISL_1068198	EPI_ISL_1068248	EPI_ISL_1068266	EPI_ISL_2777409
EPI_ISL_1068221	EPI_ISL_1068249	EPI_ISL_1068268	EPI_ISL_2777410
EPI_ISL_1068222	EPI_ISL_1068258	EPI_ISL_1068269	EPI_ISL_2777414
EPI_ISL_1068225	EPI_ISL_1068260	EPI_ISL_1068270	EPI_ISL_2777415
EPI_ISL_1068226	EPI_ISL_1068261	EPI_ISL_1068271	EPI_ISL_2777465
EPI_ISL_1068243	EPI_ISL_1068262	EPI_ISL_1068272	EPI_ISL_2777466
EPI_ISL_1068248	EPI_ISL_1068263	EPI_ISL_1068273	EPI_ISL_2777467
EPI_ISL_1068249	EPI_ISL_1068264	EPI_ISL_1068274	EPI_ISL_2777469
EPI_ISL_1068258	EPI_ISL_1068266	EPI_ISL_1068279	EPI_ISL_2777470
EPI_ISL_1068260	EPI_ISL_1068268	EPI_ISL_1068282	EPI_ISL_2777472
EPI_ISL_1068261	EPI_ISL_1068269	EPI_ISL_1068283	EPI_ISL_2777473
EPI_ISL_1068262	EPI_ISL_1068270	EPI_ISL_1068284	EPI_ISL_2777474
EPI_ISL_1068263	EPI_ISL_1068271	EPI_ISL_1068285	EPI_ISL_2777475
EPI_ISL_1068264	EPI_ISL_1068272	EPI_ISL_1068286	EPI_ISL_2777482
EPI_ISL_1068266	EPI_ISL_1068273	EPI_ISL_1068287	EPI_ISL_2777483
EPI_ISL_1068268	EPI_ISL_1068274	EPI_ISL_1068288	EPI_ISL_2777485
EPI_ISL_1068269	EPI_ISL_1068275	EPI_ISL_1068290	EPI_ISL_2777503
EPI_ISL_1068270	EPI_ISL_1068276	EPI_ISL_1068291	EPI_ISL_2777508
EPI_ISL_1068271	EPI_ISL_1068278	EPI_ISL_1068292	EPI_ISL_2777509
EPI_ISL_1068272	EPI_ISL_1068279	EPI_ISL_1166615	EPI_ISL_2777516
EPI_ISL_1068273	EPI_ISL_1068280	EPI_ISL_1213190	EPI_ISL_2777599
EPI_ISL_1068274	EPI_ISL_1068281	EPI_ISL_1213204	EPI_ISL_2777698
EPI_ISL_1068275	EPI_ISL_1068282	EPI_ISL_1261683	EPI_ISL_2777986
EPI_ISL_1068276	EPI_ISL_1068283	EPI_ISL_1261685	EPI_ISL_2777987
EPI_ISL_1068278	EPI_ISL_1068284	EPI_ISL_1261690	EPI_ISL_2777993
EPI_ISL_1068279	EPI_ISL_1068285	EPI_ISL_1261694	EPI_ISL_2777999
EPI_ISL_1068280	EPI_ISL_1068286	EPI_ISL_2777236	EPI_ISL_2778002
EPI_ISL_1068281	EPI_ISL_1068287	EPI_ISL_2777248	EPI_ISL_2778004

EPI_ISL_1068282	EPI_ISL_1068288	EPI_ISL_2777249	EPI_ISL_2778005
EPI_ISL_1068283	EPI_ISL_1068289	EPI_ISL_2777250	EPI_ISL_833138
EPI_ISL_1068284	EPI_ISL_1068290	EPI_ISL_2777320	EPI_ISL_833140
EPI_ISL_1068285	EPI_ISL_1068291	EPI_ISL_2777363	EPI_ISL_906071
EPI_ISL_1068286	EPI_ISL_1068292	EPI_ISL_2777364	EPI_ISL_918505
EPI_ISL_1068287	EPI_ISL_1166615	EPI_ISL_2777373	EPI_ISL_918506
EPI_ISL_1068288	EPI_ISL_1213190	EPI_ISL_2777374	EPI_ISL_918508
EPI_ISL_1068289	EPI_ISL_1213204	EPI_ISL_2777375	EPI_ISL_918509
EPI_ISL_1068290	EPI_ISL_1261683	EPI_ISL_2777376	
EPI_ISL_1068291	EPI_ISL_1261685	EPI_ISL_2777377	
EPI_ISL_1068292	EPI_ISL_1261690	EPI_ISL_2777378	
EPI_ISL_1166615	EPI_ISL_1261694	EPI_ISL_2777380	
EPI_ISL_1213190	EPI_ISL_2777236	EPI_ISL_2777383	
EPI_ISL_1213204	EPI_ISL_2777238	EPI_ISL_2777384	
EPI_ISL_1261683	EPI_ISL_2777248	EPI_ISL_2777385	
EPI_ISL_1261685	EPI_ISL_2777249	EPI_ISL_2777388	
EPI_ISL_1261690	EPI_ISL_2777250	EPI_ISL_2777397	
EPI_ISL_1261694	EPI_ISL_2777251	EPI_ISL_2777398	
EPI_ISL_2777236	EPI_ISL_2777320	EPI_ISL_2777399	
EPI_ISL_2777238	EPI_ISL_2777363	EPI_ISL_2777400	
EPI_ISL_2777248	EPI_ISL_2777364	EPI_ISL_2777401	
EPI_ISL_2777249	EPI_ISL_2777373	EPI_ISL_2777402	
EPI_ISL_2777250	EPI_ISL_2777374	EPI_ISL_2777403	
EPI_ISL_2777251	EPI_ISL_2777375	EPI_ISL_2777404	
EPI_ISL_2777320	EPI_ISL_2777376	EPI_ISL_2777405	
EPI_ISL_2777363	EPI_ISL_2777377	EPI_ISL_2777406	
EPI_ISL_2777364	EPI_ISL_2777378	EPI_ISL_2777407	
EPI_ISL_2777373	EPI_ISL_2777380	EPI_ISL_2777408	
EPI_ISL_2777374	EPI_ISL_2777382	EPI_ISL_2777410	
EPI_ISL_2777375	EPI_ISL_2777383	EPI_ISL_2777412	
EPI_ISL_2777376	EPI_ISL_2777384	EPI_ISL_2777413	
EPI_ISL_2777377	EPI_ISL_2777385	EPI_ISL_2777414	

EPI_ISL_2777378	EPI_ISL_2777388	EPI_ISL_2777415	
EPI_ISL_2777380	EPI_ISL_2777397	EPI_ISL_2777417	
EPI_ISL_2777382	EPI_ISL_2777398	EPI_ISL_2777418	
EPI_ISL_2777383	EPI_ISL_2777399	EPI_ISL_2777419	
EPI_ISL_2777384	EPI_ISL_2777400	EPI_ISL_2777454	
EPI_ISL_2777385	EPI_ISL_2777401	EPI_ISL_2777461	
EPI_ISL_2777388	EPI_ISL_2777402	EPI_ISL_2777462	
EPI_ISL_2777397	EPI_ISL_2777403	EPI_ISL_2777465	
EPI_ISL_2777398	EPI_ISL_2777404	EPI_ISL_2777466	
EPI_ISL_2777399	EPI_ISL_2777405	EPI_ISL_2777467	
EPI_ISL_2777400	EPI_ISL_2777406	EPI_ISL_2777469	
EPI_ISL_2777401	EPI_ISL_2777407	EPI_ISL_2777470	
EPI_ISL_2777402	EPI_ISL_2777408	EPI_ISL_2777472	
EPI_ISL_2777403	EPI_ISL_2777409	EPI_ISL_2777473	
EPI_ISL_2777404	EPI_ISL_2777410	EPI_ISL_2777474	
EPI_ISL_2777405	EPI_ISL_2777412	EPI_ISL_2777475	
EPI_ISL_2777406	EPI_ISL_2777413	EPI_ISL_2777477	
EPI_ISL_2777407	EPI_ISL_2777414	EPI_ISL_2777478	
EPI_ISL_2777408	EPI_ISL_2777415	EPI_ISL_2777479	
EPI_ISL_2777409	EPI_ISL_2777416	EPI_ISL_2777481	
EPI_ISL_2777410	EPI_ISL_2777417	EPI_ISL_2777482	
EPI_ISL_2777412	EPI_ISL_2777418	EPI_ISL_2777483	
EPI_ISL_2777413	EPI_ISL_2777419	EPI_ISL_2777485	
EPI_ISL_2777414	EPI_ISL_2777420	EPI_ISL_2777495	
EPI_ISL_2777415	EPI_ISL_2777454	EPI_ISL_2777498	
EPI_ISL_2777416	EPI_ISL_2777460	EPI_ISL_2777503	
EPI_ISL_2777417	EPI_ISL_2777461	EPI_ISL_2777507	
EPI_ISL_2777418	EPI_ISL_2777462	EPI_ISL_2777508	
EPI_ISL_2777419	EPI_ISL_2777464	EPI_ISL_2777539	
EPI_ISL_2777420	EPI_ISL_2777466	EPI_ISL_2777599	
EPI_ISL_2777454	EPI_ISL_2777467	EPI_ISL_2777698	
EPI_ISL_2777460	EPI_ISL_2777468	EPI_ISL_2777700	

EPI_ISL_ 2777461	EPI_ISL_ 2777469	EPI_ISL_ 2777701	
EPI_ISL_ 2777462	EPI_ISL_ 2777470	EPI_ISL_ 2777740	
EPI_ISL_ 2777464	EPI_ISL_ 2777472	EPI_ISL_ 2777986	
EPI_ISL_ 2777465	EPI_ISL_ 2777473	EPI_ISL_ 2777987	
EPI_ISL_ 2777466	EPI_ISL_ 2777475	EPI_ISL_ 2777993	
EPI_ISL_ 2777467	EPI_ISL_ 2777477	EPI_ISL_ 2777995	
EPI_ISL_ 2777468	EPI_ISL_ 2777478	EPI_ISL_ 2777996	
EPI_ISL_ 2777469	EPI_ISL_ 2777481	EPI_ISL_ 2777997	
EPI_ISL_ 2777470	EPI_ISL_ 2777482	EPI_ISL_ 2777998	
EPI_ISL_ 2777471	EPI_ISL_ 2777495	EPI_ISL_ 2777999	
EPI_ISL_ 2777472	EPI_ISL_ 2777498	EPI_ISL_ 2778000	
EPI_ISL_ 2777473	EPI_ISL_ 2777499	EPI_ISL_ 2778002	
EPI_ISL_ 2777474	EPI_ISL_ 2777503	EPI_ISL_ 2778005	
EPI_ISL_ 2777475	EPI_ISL_ 2777508	EPI_ISL_ 811149	
EPI_ISL_ 2777477	EPI_ISL_ 2777516	EPI_ISL_ 833136	
EPI_ISL_ 2777478	EPI_ISL_ 2777539	EPI_ISL_ 833139	
EPI_ISL_ 2777479	EPI_ISL_ 2777698	EPI_ISL_ 833140	
EPI_ISL_ 2777481	EPI_ISL_ 2777701	EPI_ISL_ 906071	
EPI_ISL_ 2777482	EPI_ISL_ 2777740	EPI_ISL_ 906077	
EPI_ISL_ 2777483	EPI_ISL_ 2777986	EPI_ISL_ 906081	
EPI_ISL_ 2777484	EPI_ISL_ 2777987	EPI_ISL_ 918500	
EPI_ISL_ 2777485	EPI_ISL_ 2777995	EPI_ISL_ 918502	
EPI_ISL_ 2777495	EPI_ISL_ 2777996	EPI_ISL_ 918503	
EPI_ISL_ 2777498	EPI_ISL_ 2777997	EPI_ISL_ 918506	
EPI_ISL_ 2777499	EPI_ISL_ 2777998	EPI_ISL_ 918508	
EPI_ISL_ 2777503	EPI_ISL_ 2778002	EPI_ISL_ 918509	
EPI_ISL_ 2777507	EPI_ISL_ 2778005	EPI_ISL_ 918511	
EPI_ISL_ 2777508	EPI_ISL_ 811149		
EPI_ISL_ 2777509	EPI_ISL_ 833136		
EPI_ISL_ 2777516	EPI_ISL_ 833138		
EPI_ISL_ 2777539	EPI_ISL_ 833139		
EPI_ISL_ 2777599	EPI_ISL_ 833140		



EPI_ISL_ 2777698	EPI_ISL_ 906071		
EPI_ISL_ 2777700	EPI_ISL_ 906080		
EPI_ISL_ 2777701	EPI_ISL_ 906081		
EPI_ISL_ 2777740	EPI_ISL_ 918500		
EPI_ISL_ 2777986	EPI_ISL_ 918501		
EPI_ISL_ 2777987	EPI_ISL_ 918502		
EPI_ISL_ 2777993	EPI_ISL_ 918503		
EPI_ISL_ 2777995	EPI_ISL_ 918505		
EPI_ISL_ 2777996	EPI_ISL_ 918506		
EPI_ISL_ 2777997	EPI_ISL_ 918507		
EPI_ISL_ 2777998	EPI_ISL_ 918508		
EPI_ISL_ 2777999	EPI_ISL_ 918510		
EPI_ISL_ 2778000	EPI_ISL_ 918511		
EPI_ISL_ 2778002			
EPI_ISL_ 2778004			
EPI_ISL_ 2778005			
EPI_ISL_ 811149			
EPI_ISL_ 833136			
EPI_ISL_ 833138			
EPI_ISL_ 833139			
EPI_ISL_ 833140			
EPI_ISL_ 906071			
EPI_ISL_ 906075			
EPI_ISL_ 906076			
EPI_ISL_ 906077			
EPI_ISL_ 906080			
EPI_ISL_ 906081			
EPI_ISL_ 918499			
EPI_ISL_ 918500			
EPI_ISL_ 918501			
EPI_ISL_ 918502			
EPI_ISL_ 918503			

EPI_ISL_ 918504			
EPI_ISL_ 918505			
EPI_ISL_ 918506			
EPI_ISL_ 918507			
EPI_ISL_ 918508			
EPI_ISL_ 918509			
EPI_ISL_ 918510			
EPI_ISL_ 918511			