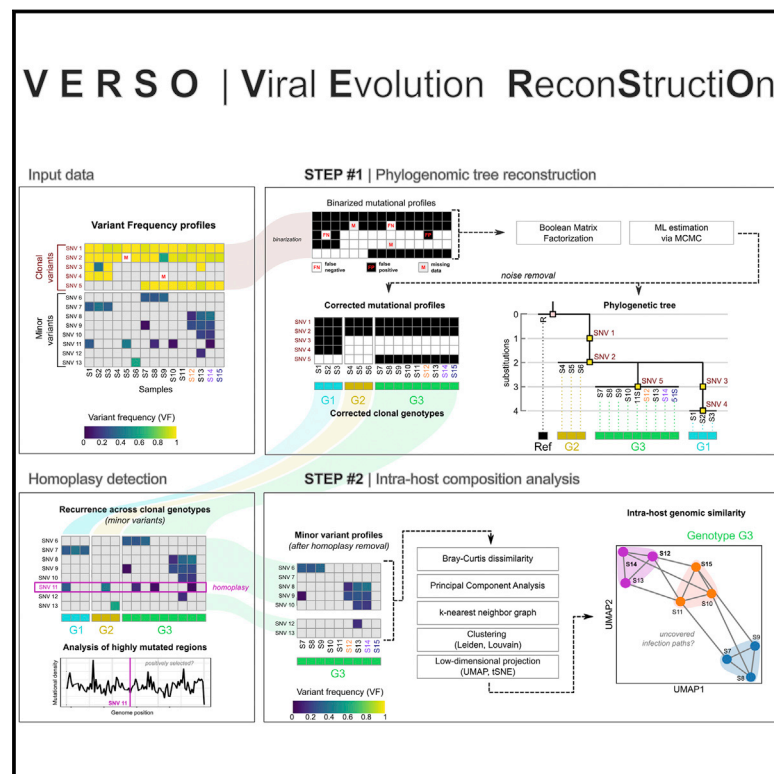


Patterns

VERSO: A comprehensive framework for the inference of robust phylogenies and the quantification of intra-host genomic diversity of viral samples

Graphical abstract



Highlights

- The analysis of raw sequencing data improves the reconstruction of viral evolution
- Our method reconstructs robust phylogenies with noisy data and sampling limitations
- The dissection of intra-host genomic diversity reveals undetected infection chains
- The identification of positively selected variants may drive experimental research

Authors

Daniele Ramazzotti, Fabrizio Angaroni, Davide Maspero, Carlo Gambacorti-Passerini, Marco Antoniotti, Alex Graudenzi, Rocco Piazza

Correspondence

alex.graudenzi@ibfm.cnr.it (A.G.), rocco.piazza@unimib.it (R.P.)

In Brief

The generation of reliable phylogenomic models describing the evolution of SARS-CoV-2 is essential to explain its diffusion and to possibly predict the next evolutionary steps. We introduce a data-science framework that is an improvement on existing methods, by accounting for noise and sampling limitations in sequencing data and by dissecting the intra-host diversity of single samples. The application to large-scale datasets demonstrates that our approach can improve the estimation of SARS-CoV-2 evolution, refine contact tracing, and pinpoint possibly hazardous mutations.



Article

VERSO: A comprehensive framework for the inference of robust phylogenies and the quantification of intra-host genomic diversity of viral samples

Daniele Ramazzotti,¹ Fabrizio Angaroni,² Davide Maspero,^{2,3} Carlo Gambacorti-Passerini,¹ Marco Antoniotti,^{2,4} Alex Graudenzi,^{3,4,5,6,*} and Rocco Piazza^{1,5,*}

¹Department of Medicine and Surgery, Università degli Studi di Milano-Bicocca, Monza, Italy

²Department of Informatics, Systems and Communication, Università degli Studi di Milano-Bicocca, Milan, Italy

³Inst. of Molecular Bioimaging and Physiology, Consiglio Nazionale delle Ricerche (IBFM-CNR), Segrate, Milan, Italy

⁴Bicocca Bioinformatics, Biostatistics and Bioimaging Centre – B4, Milan, Italy

⁵Senior author

⁶Lead contact

*Correspondence: alex.graudenzi@ibfm.cnr.it (A.G.), rocco.piazza@unimib.it (R.P.)

<https://doi.org/10.1016/j.patter.2021.100212>

THE BIGGER PICTURE The gravity of the COVID-19 pandemic has fostered a surge of works analyzing SARS-CoV-2 consensus sequences to reconstruct phylogenomic models of its evolution and diffusion. Yet, such approaches do not account for intra-host genomic diversity and may deliver inaccurate predictions in conditions of noisy data and sampling limitations.

We propose VERSO, a data-science framework for the characterization of viral evolution from sequencing data. By accounting for uncertainty in the data, VERSO delivers robust phylogenies also in conditions of limited sampling and noisy observations. Additionally, the in-depth characterization of the intra-host genomic diversity of samples allows one to identify undetected infection chains and clusters and to intercept variants possibly undergoing positive selection. Accordingly, the joint application of our method and data-driven epidemiological models may deliver a high-precision platform for contact tracing and pathogen surveillance and characterization.



Proof-of-Concept: Data science output has been formulated, implemented, and tested for one domain/problem

SUMMARY

We introduce VERSO, a two-step framework for the characterization of viral evolution from sequencing data of viral genomes, which is an improvement on phylogenomic approaches for consensus sequences. VERSO exploits an efficient algorithmic strategy to return robust phylogenies from clonal variant profiles, also in conditions of sampling limitations. It then leverages variant frequency patterns to characterize the intra-host genomic diversity of samples, revealing undetected infection chains and pinpointing variants likely involved in homoplasies. On simulations, VERSO outperforms state-of-the-art tools for phylogenetic inference. Notably, the application to 6,726 amplicon and RNA sequencing samples refines the estimation of severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2) evolution, while co-occurrence patterns of minor variants unveil undetected infection paths, which are validated with contact tracing data. Finally, the analysis of SARS-CoV-2 mutational landscape uncovers a temporal increase of overall genomic diversity and highlights variants transiting from minor to clonal state and homoplastic variants, some of which fall on the spike gene. Available at: <https://github.com/BIMIB-DISCO/VERSO>.

INTRODUCTION

The outbreak of coronavirus disease 2019 (COVID-19), which started in late 2019 in Wuhan (China)^{1,2} and was declared a

pandemic by the World Health Organization, is fueling the publication of an increasing number of studies aimed at exploiting the information provided by the viral genome of severe acute respiratory syndrome-coronavirus 2 (SARS-CoV-2) virus to identify its



proximal origin, characterize the mode and timing of its evolution, as well as to define descriptive and predictive models of geographical spread and evaluate the related clinical impact.^{3–5} As a matter of fact, the mutations that rapidly accumulate in the viral genome⁶ can be used to track the evolution of the virus and, accordingly, unravel the viral infection network.^{7,8}

At the time of this writing, numerous independent laboratories around the world are isolating and sequencing SARS-CoV-2 samples and depositing them on public databases (e.g., GISAID⁹) whose data are accessible via the Nextstrain portal.¹⁰ Such data can be employed to estimate models from genomic epidemiology and may serve, for instance, to estimate the proportion of undetected infected people by uncovering cryptic transmissions, as well as to predict likely trends in the number of infected, hospitalized, dead, and recovered people.^{11–13}

More in detail, most studies employ phylogenomic approaches that process consensus sequences, which represent the dominant virus lineage within each infected host. A growing plethora of methods for phylogenomic reconstruction is available to this end, all relying on different algorithmic frameworks, including distance-matrix, maximum parsimony, maximum likelihood, or Bayesian inference, with various substitution models and distinct evolutionary assumptions (see, e.g., Refs.^{10,14–22}). However, while such methods have repeatedly proven effective in unraveling the main patterns of evolution of viral genomes with respect to many different diseases, including SARS-CoV-2,^{10,23–25} at least two issues can be raised.

First, most phylogenomics methods might produce unreliable results when dealing with noisy data, for instance due to sequencing issues, or with data collected with significant sampling limitations,^{14,26,27} as witnessed for most countries during the epidemics.^{28,29}

Second, most methods do not consider the key information on intra-host minor variants (also referred to as minority variants or intra-host single nucleotide variants), which can be retrieved from whole-genome deep sequencing raw data and might be essential to improve the characterization of the infection dynamics and to pinpoint positively selected variants.^{30–32} Due to the high replication, mutation, and recombination rates of RNA viruses, subpopulations of mutant viruses, also known as viral quasispecies,³⁰ typically emerge and coexist within single hosts, and are supposed to underlie most of the adaptive potential to the immune system response and to anti-viral therapies.^{31,33,34} In this regard, many recent studies highlighted the noteworthy amount of intra-host genomic diversity in SARS-CoV-2 samples,^{35–43} similarly to what has already been observed in many distinct infectious diseases.^{8,32,44–48}

Here, we introduce VERSO (viral evolution reconstruction), a new comprehensive framework for the inference of high-resolution models of viral evolution from raw sequencing data of viral genomes (see Figure 1). VERSO includes two consecutive algorithmic steps.

Step #1: robust phylogenomic inference from clonal variant profiles

VERSO first employs a probabilistic noise-tolerant framework to process binarized clonal variant profiles (or, alternatively, consensus sequences), to return a robust phylogenetic model also in conditions of sampling limitations and sequencing issues.

By adapting algorithmic strategies widely employed in cancer evolution analysis,^{49–52} VERSO is able to correct false-positive and false-negative variants, can manage missing observations due to low coverage, and is designed to group samples with identical (corrected) clonal genotype in polytomies, avoiding ungrounded arbitrary orderings. As a result, the accurate and robust phylogenomic models produced by VERSO may be used to improve the parameter estimation of epidemiological models, which typically rely on limited and inhomogeneous data.^{11,29} Notice that this step can be executed independently from step #2; for instance, in case raw sequencing data are not available.

Homoplasmy detection (clonal variants)

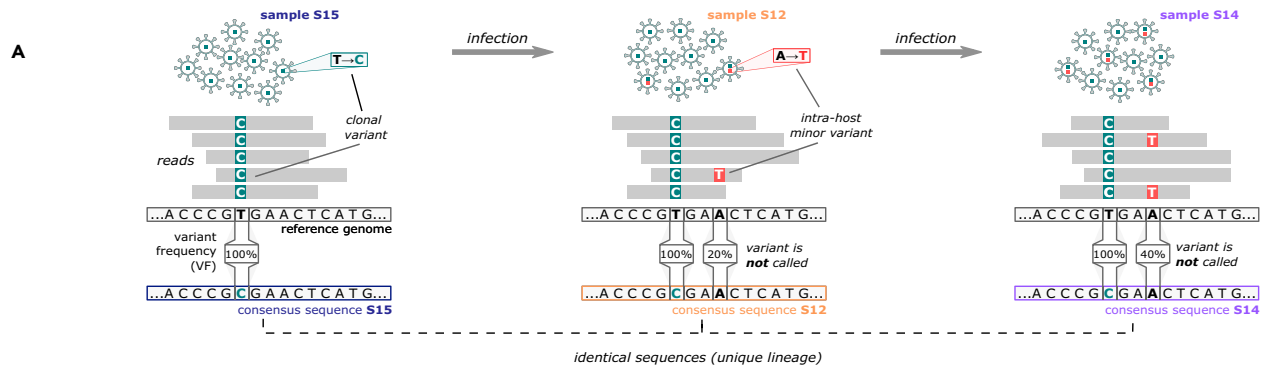
The first step of VERSO allows one to identify clonal mutations that might be involved in reticulation events^{53,54} and, in particular, in homoplasies, possibly due to positive selection in a scenario of convergent/parallel evolution,⁵⁵ founder effects,³¹ or mutational hotspots.⁵⁶ Such information might be useful to drive the design of opportune treatments and vaccines; for instance, by blacklisting positively selected genomic regions.

Step #2: characterization of intra-host genomic diversity

In the second step, VERSO exploits the information on variant frequency (VF) profiles obtained from raw sequencing data (if available), to characterize and visualize the intra-host genomic similarity of hosts with identical (corrected) clonal genotype. In fact, even though the extent and modes of transmission of quasispecies from a host to another during infections are still elusive,^{31,57} patterns of co-occurrence of minor variants detected in hosts with identical clonal genotype may provide an indication on the presence of undetected infection paths.^{8,58} For this reason, the second step of VERSO is designed to characterize and visualize the genomic similarity of samples by exploiting dimensionality reduction and clustering strategies typically employed in single-cell analyses.⁵⁹ Alternative approaches for the analysis of quasispecies, yet with different goals and algorithmic assumptions, have been proposed, for instance in Refs.^{60–63} and recently reviewed in Knyazev et al.⁶⁴ As specified above, VERSO step #2 is executed on groups of samples with identical clonal genotype: the rationale is that the transmission of minor variants implicates the concurrent transfer of clonal variants, excluding the rare cases in which the VF of a clonal variant significantly decreases in a given host; for instance due to mutation losses (e.g., via recombination-associated deletions or via multiple mutations hitting an already mutated genome location³⁴) or to complex horizontal evolution phenomena (e.g., super-infections^{65,66}). Conversely, the transmission of clonal variants does not necessarily implicate the transfer of all minor variants, which are affected by complex recombination and transmission effects, such as bottlenecks.^{31,57} As a final result, VERSO allows one to visualize the genomic similarity of samples on a low-dimensional space (e.g., UMAP [uniform manifold approximation and projection]⁶⁷ or tSNE [t-distributed stochastic neighbor embedding]⁶⁸) representing the intra-host genomic diversity, and to characterize high-resolution infection chains, thus overcoming the limitations of methods relying on consensus sequences.

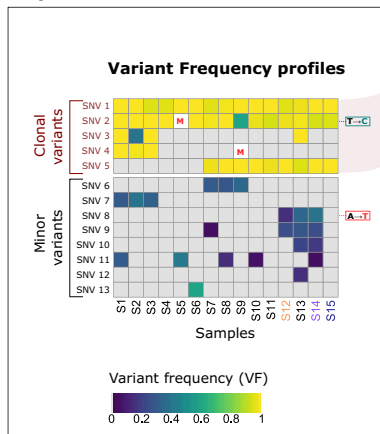
Homoplasmy detection (minor variants)

Importantly, minor variants observed in hosts with distinct clonal genotypes (identified via VERSO step #1) may indicate

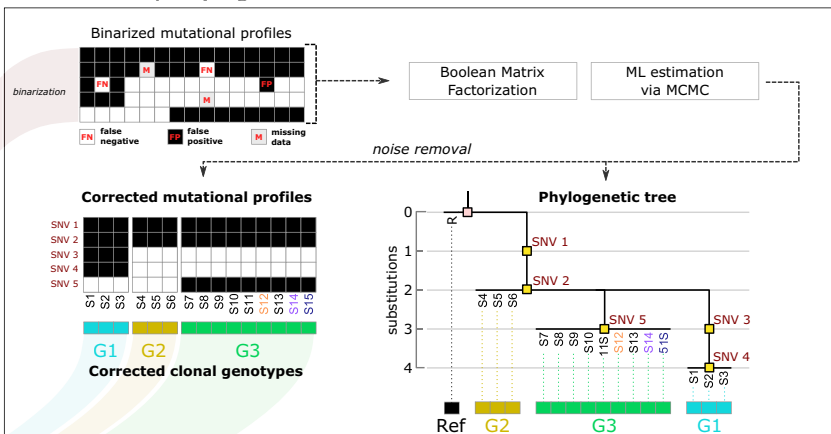


B VERSO | Viral Evolution ReconStruction

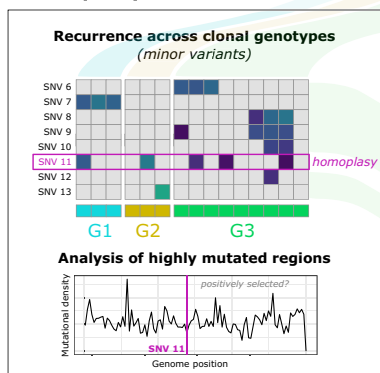
Input data



STEP #1 | Phylogenomic tree reconstruction



Homoplasy detection



STEP #2 | Intra-host composition analysis

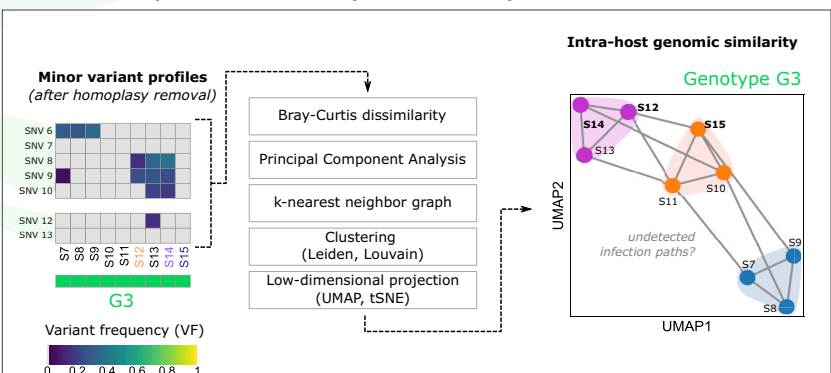
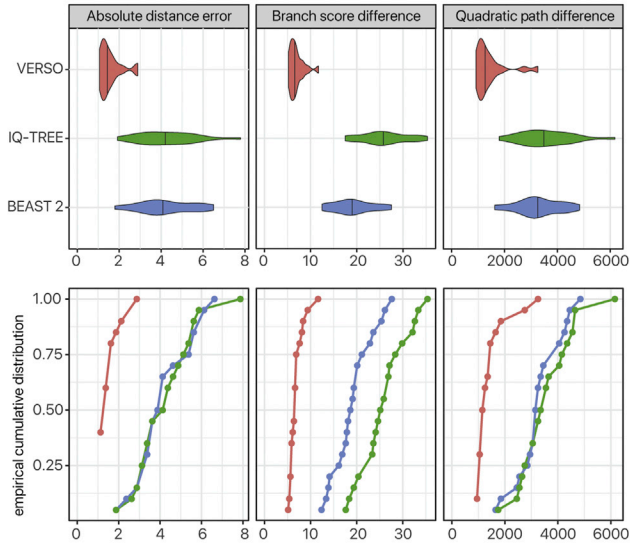


Figure 1. VERSO framework for viral evolution inference and intra-host genomic diversity quantification

(A) In this example, three hosts infected by the same viral lineage are sequenced. All hosts share the same clonal mutation (T>C, green), but two of them (#2 and #3) are characterized by a distinct minor mutation (A>T, red), which randomly emerged in host #2 and was transferred to host #3 during the infection. Standard sequencing experiments return an identical consensus sequence for all samples, by employing a threshold on VF and by selecting mutations characterizing the dominant lineage.

(B) VERSO takes as input the VF profiles of samples, generated from raw sequencing data. In step #1, VERSO processes the binarized profiles of clonal variants and solves a Boolean matrix factorization problem by maximizing a likelihood function via MCMC, in order to correct false-positives/-negatives and missing data. As output, it returns both the corrected mutational profiles of samples and the phylogenetic tree, in which samples with identical corrected clonal genotypes are grouped in polytomies. Corrected clonal genotypes are then employed to identify homoplasies of minor variants, which are further investigated to pinpoint positively selected mutations. The VF profile of minor variants (excluding homoplasies) is processed by step #2 of VERSO, which computes a refined genomic distance among hosts (via Bray-Curtis dissimilarity, after PCA) and performs clustering and dimensionality reduction, in order to project and visualize samples on a 2D space, representing the intra-host genomic diversity and the distance among hosts. This allows one to identify undetected transmission paths among samples with identical clonal genotype.

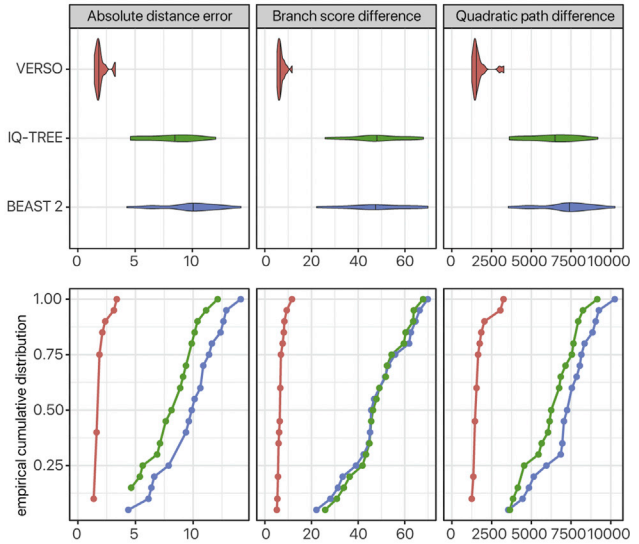
A Low noise - No sampling limitations



	Improv. median	U-test 2-sided p-Value	Improv. median	U-test 2-sided p-Value	Improv. median	U-test 2-sided p-Value
IQ-TREE	-67.7%	< .001	-74.3%	< .001	-65.4%	< .001
BEAST 2	-66.2%	< .001	-65.3%	< .001	-62.4%	< .001

20 topologies: n = 1000, $\alpha = 0.05$, $\beta = 0.05$, $m \in \{14, \dots, 31\}$

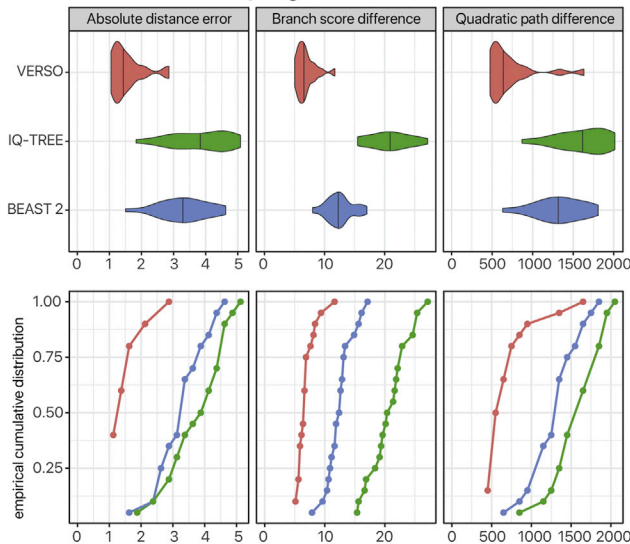
B High noise - No sampling limitations



	Improv. median	U-test 2-sided p-Value	Improv. median	U-test 2-sided p-Value	Improv. median	U-test 2-sided p-Value
IQ-TREE	-78.6%	< .001	-86.1%	< .001	-76.6%	< .001
BEAST 2	-81.7%	< .001	-85.8%	< .001	-79.5%	< .001

20 topologies: n = 1000, $\alpha = 0.10$, $\beta = 0.10$, $m \in \{14, \dots, 31\}$

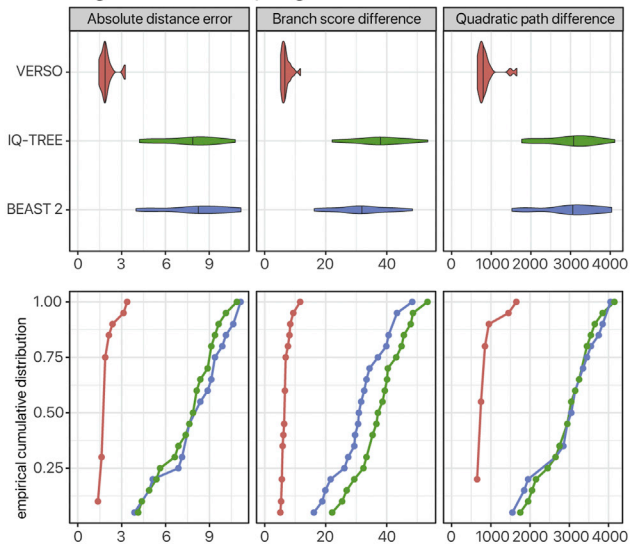
C Low noise - Sampling limitations



	Improv. median	U-test 2-sided p-Value	Improv. median	U-test 2-sided p-Value	Improv. median	U-test 2-sided p-Value
IQ-TREE	-66.4%	< .001	-68.5%	< .001	-64.9%	< .001
BEAST 2	-59.8%	< .001	-47.3%	< .001	-56.3%	< .001

20 topologies: n = 500, $\alpha = 0.05$, $\beta = 0.05$, $m \in \{14, \dots, 31\}$

D High noise - Sampling limitations



	Improv. median	U-test 2-sided p-Value	Improv. median	U-test 2-sided p-Value	Improv. median	U-test 2-sided p-Value
IQ-TREE	-76.6%	< .001	-82.7%	< .001	-74.9%	< .001
BEAST 2	-78.1%	< .001	-79.0%	< .001	-75.3%	< .001

20 topologies: n = 500, $\alpha = 0.10$, $\beta = 0.10$, $m \in \{14, \dots, 31\}$

■ VERSO ■ IQ-TREE ■ BEAST 2

Figure 2. Comparative assessment on simulated data

(A–D) Synthetic datasets were generated via the widely used coalescent model simulator msprime⁷⁰ (see the Supplementary Material and Table S1 for the parameter settings). Twenty distinct topologies with 1,000 samples were generated, including a number of distinguishable variants in the range (14, 31). For each topology, four synthetic datasets were generated, with different sample sizes (n = 1000, 500), and different combinations of false-positives and false-negatives

(legend continued on next page)

homoplasies, due to mutational hotspots, phantom mutations, or to positive selection.⁵⁶ VERSO pinpoints such variants for further investigations and allows one to exclude them from the computation of the VF-based genomic similarity prior to VERSO step #2, to reduce the possible confounding effects.

To summarize, VERSO (1) returns accurate and robust phylogenies of viral samples, by removing noise from clonal variant profiles; (2) detects reticulation events due to homoplasies of clonal variants; (3) exploits minor variant profiles to characterize and visualize the intra-host genomic similarity of samples with identical (corrected) clonal genotype, thus pinpointing undetected infection paths; (4) allows one to identify and characterize homoplastic minor variants, which might be due to positive selection or mutational hotspots.

To assess the accuracy and robustness of the results produced by VERSO, we performed an extensive array of simulations, and compared with two state-of-the-art methods for phylogenetic reconstruction; i.e., IQ-TREE¹⁰ and BEAST 2.²² As a major result, VERSO outperforms competing methods in all settings and also in condition of high noise and sampling limitations.

Furthermore, we applied VERSO to two large-scale datasets, generated via amplicon and RNA-seq Illumina sequencing protocols, including 3,960 and 2,766 samples, respectively. The robust phylogenomic models delivered via VERSO step #1 allow us to refine the current estimation on SARS-CoV-2 evolution and spread. Besides, thanks to the in-depth analysis of the mutational landscape of both clonal and minor variants, we could identify a number of variants undergoing transition to clonality, as well as several homoplasies, including variants likely undergoing positive selection processes.

Remarkably, the infection chains identified via VERSO step #2, by assessing the intra-host genomic similarity of samples with the same clonal genotype, were validated by employing contact tracing data from Rockett et al.⁶⁹ This important result, which could not be achieved by analyzing consensus sequences, proves the effectiveness of employing raw sequencing data to improve the characterization of the transmission dynamics, in particular during the early phase of the outbreak, in which a relatively low diversity of SARS-CoV-2 has been observed at the consensus level.

VERSO is released as free open source tool at this link: <https://github.com/BIMIB-DISCO/VERSO>.

RESULTS

Comparative assessment on simulations

In order to assess the performance of VERSO and compare it with competing approaches, we executed extensive tests on simulated datasets, generated with the coalescent model simulator *msprime*.⁷⁰ Simulations allow one to compute a number of metrics with respect to the ground truth, which in this case is the phylogeny of samples resulting from a backwards-in-time coalescent simulation.⁷¹ Accordingly, this allows one to evaluate

the accuracy and robustness of the results produced by competing methods in a variety of in-silico scenarios.

In detail, we selected 20 simulation scenarios with $n = 1,000$ samples in which a number of clonal variants (with distinguishable profiles) between 14 and 31 was observed. We then inflated the datasets with false-positives with rate α and false-negatives with rate β , in order to mimic sequencing and coverage issues. Moreover, additional datasets were generated via random subsampling of the original datasets, to model possible sampling limitations and sampling biases. As a result, we investigated four simulation settings: (A) low noise, no subsampling; (B) high noise, no subsampling; (C) low noise, subsampling; and (D) high noise, subsampling (see [Experimental procedures](#) and the [Supplemental experimental procedures](#) for further details; the complete parameter settings of the simulations are provided in [Table S1](#)).

VERSO step #1 was compared with two state-of-the-art phylogenetic methods from consensus sequences: IQ-TREE,¹⁰ the algorithmic strategy included in the Nextstrain-Augur pipeline,⁷² and BEAST 2.²² Consensus sequences to be provided as input to such methods were generated from simulation data by employing the reference genome SARS-CoV-2-ANC (see below).

The performance of methods was assessed by comparing the reconstructed phylogeny with the simulated ground truth, in terms of (1) absolute error evolutionary distance, (2) branch score difference,⁷³ and (3) quadratic path difference⁷⁴ (please refer to the [Supplemental experimental procedures](#) for a detailed description of all metrics).

[Figure 2](#) shows the performance distribution of all methods with respect to all simulation settings. Notably, VERSO step #1 outperforms competing methods in all scenarios (Mann-Whitney U test, $p < 0.001$ in all cases), with noteworthy percentage improvements, also in conditions of high noise and sampling limitations. This important result shows that the probabilistic framework that underlies VERSO step #1 can produce more robust and reliable results when processing noisy data, as typically observed in real-world scenarios.

Reference genome

Different reference genomes have been employed in the analysis of SARS-CoV-2 origin and evolution. Two genome sequences from human samples, in particular, were used in early phylogenomic studies, namely sequence EPI_ISL_405839 (ref. #1 in the following) used, e.g., in Bastola et al.⁷⁵ and sequence EPI_ISL_402125 (ref. #2) used, e.g., in Andersen et al.³ Excluding the polyA tails, the two sequences are identical for 29,865 of 29,870 genome positions (99.98%) and differ for only five SNPs at locations 8,782, 9,561, 15,607, 28,144, and 29,095, for which ref. #1 has haplotype TTCCT and ref. #2 has haplotype CCTTC.

In order to define a likely common ancestor for both sequences, we analyzed the Bat-CoV-RaTG13 genome (sequence

($[\alpha = 0.05, \beta = 0.05]$, $[\alpha = 0.10, \beta = 0.10]$), for a total of four configurations (A, B, C, and D) and 80 independent datasets. VERSO step #1 was compared with IQ-TREE¹⁰ and BEAST 2,²² on (1) absolute error evolutionary distance, (2) branch score difference⁷³ and (3) quadratic path difference⁷⁴ with respect to the ground-truth sample phylogeny provided by *msprime* (see the [Supplemental experimental procedures](#) for the description of the metrics). In the upper panels, distributions are shown as violin plots, whereas lower panels include the empirical cumulative distribution functions. The percentage improvement of VERSO with respect to competing methods is shown on all metrics (computed on median values), in addition to the p value of the two-sided Mann-Whitney U test on distributions, for all settings.

EPI_ISL_402131)¹ and the Pangolin-CoV genome (sequence EPI_ISL_410721),^{3,4} which were identified as closely related genomes to SARS-CoV-2.⁷⁶ In particular, it was hypothesized that SARS-CoV-2 might be a recombinant of an ancestor of Pangolin-CoV and Bat-CoV-RaTG13,^{4,77} whereas more recent findings would suggest that the SARS-CoV-2 lineage is the consequence of a direct or indirect zoonotic jump from bats.⁷⁶ Whatever the case, both Bat-CoV-RaTG13 and Pangolin-CoV display haplotype TCTCT at locations 8,782, 9,561, 15,607, 28,144 and 29,095 and, therefore, one can hypothesize that such a haplotype was present in the unknown common ancestor of ref. #1 and #2.

For this reason, we generated an artificial reference genome, named SARS-CoV-2-ANC, which is identical to both ref. #1 and #2 on 29,865 (over 29,870) genome locations, includes the polyA tail of ref. #2 (33 bases), and has haplotype TCTCT at locations 8,782, 9,561, 15,607, 28,144, and 29,095 (see Figure S2 for a depiction of the artificial genome generation). SARS-CoV-2-ANC is a likely common ancestor of both genomes and was used for variant calling in downstream analyses (SARS-CoV-2-ANC is released in FASTA format as Data S1). Notice that VERSO pipeline is flexible and can employ any reference genome.

Application of VERSO to 3,960 samples from amplicon sequencing data (dataset #1)

We retrieved raw Illumina Amplicon sequencing data of 3,960 SARS-CoV-2 samples of dataset #1 and applied VERSO to the mutational profiles of 2,906 samples selected after quality check (mutational profiles were generated by executing variant calling via standard practices; see Experimental procedures for further details). Notice that the analysis of this dataset was performed independently from that of dataset #2 in order to exclude possible sequencing-related artifacts or idiosyncrasies.

VERSO step #1: robust phylogenomic inference from clonal variant profiles

We first applied VERSO step #1 to the mutational profile of the 29 variants detected as clonal (VF > 90%) in at least 3% of the samples, in order to reconstruct a robust phylogenomic tree. The VERSO phylogenetic model is displayed in Figure 3A and highlights the presence of 25 clonal genotypes, obtained by removing noise from data, and that define polytomies including different numbers of samples (see Experimental procedures for further details). The mapping between clonal genotype labels and the lineage dynamic nomenclature proposed by Rambaut et al.⁷⁸ was obtained via pangolin 2.0⁷⁹ and is provided in Data S3.

More in detail, variant g.29095T>C (N, synonymous) is the earliest evolutionary event from reference genome SARS-CoV-2-ANC and is detected in 2,454 samples of the dataset. The related clonal genotype G1, which is characterized by no further mutations, identifies a polytomy including 57 Australian, 15 Chinese, 12 American, and one South-African samples.

Three clades originate from G1: a first clade includes clonal genotypes G2 (six samples) and G3 (103), while a second clade includes clonal genotype G4 (86). Clonal genotypes G1ffiG4 are characterized by the absence of single nucleotide variants (SNVs) g.8782T>C (*ORF1ab*, synonymous) and g.28144C>T (*ORF8*, p.845>L) and correspond to previously identified type

A²⁴ (also type S⁸²), which was hypothesized to be an early SARS-CoV-2 type.

The third clade originating from clonal genotype G1 includes all remaining clonal genotypes (G5–G25) and is characterized by the presence of both SNVs g.8782T>C and g.28144C>T. This specific haplotype corresponds to type B²⁴ (also type L⁸²) and an increase of its prevalence has progressively recorded in the population, as one can see in Figure 3, as opposed to type A (S), which was rarely observed in late samples. In this regard, we note that there are currently insufficient elements to support any epidemiological claim on virulence and pathogenicity of such SARS-CoV-2 types, even if recent evidences would suggest the existence of a low correlation.⁸³

Variant g.23403A>G (S, p.614D>G) is of particular interest, as proven by the increasing number of related studies.^{84–87} Such a variant identifies a large clade including 11 clonal genotypes: G15 (493 samples), G16 (1), G17 (512), G18 (25), G19 (118), G20 (94), G21 (648), G22 (90), G23 (127), G24 (4), and G25 (86), for a total of 2,198 total samples, distributed especially in Australia (971), the United States (841), South Africa (257), and Israel (125). Importantly, a constant increase of the prevalence of the haplotype corresponding to such variant is observed in time (see Figure 3), which might hint at ongoing positive selection processes; e.g., due to increased viral transmission. However, this hypothesis is highly debated⁸⁸ and, in order to investigate the possible functional effect of such variant and the related clinical implications, *in vivo* and *in vitro* studies are needed.⁸⁷

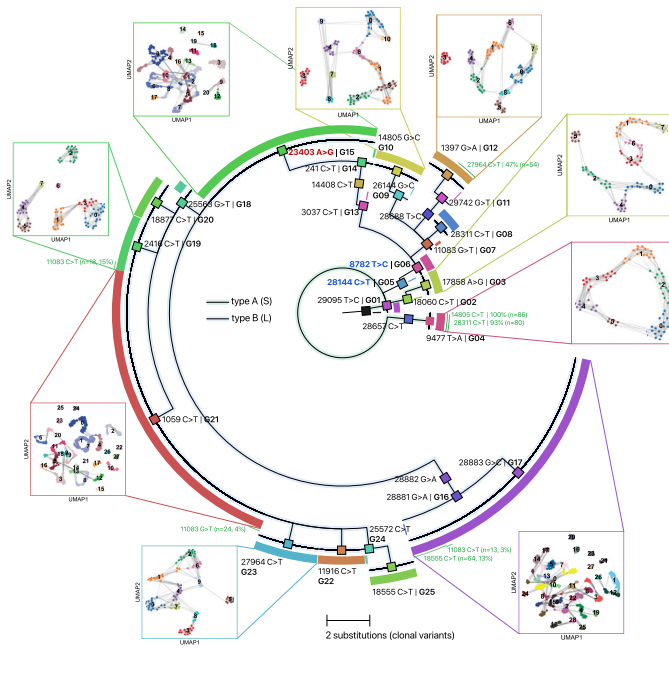
By looking more in detail at the geo-temporal localization of samples depicted via Microreact⁸¹ (Figure 3B), one can see that the different clonal genotypes are distributed across the world in distinct complex patterns, suggesting that most countries might have suffered from multiple introductions, especially in the early phases of the epidemics. In particular, samples are distributed in 11 countries, with Australia (1,523 samples), United States (910), South Africa (260), Israel (133), and China (45) representing around 99% of the dataset.

The country displaying the largest number of samples is Australia, with 1,523 samples, distributed in 22 different clonal genotypes. The presence of a number of early clonal genotypes (i.e., G1, G2, G3, G4, and G6) supports the hypothesis of multiple introductions of SARS-CoV-2 in Australia. Interestingly, we note that, from the 16th week on, the composition of the Australian sample group tends to be polarized toward clonal genotypes G17 (108/311 ≈ 35%) and G25 (82/311 ≈ 26%).

910 samples from the United States are included in the dataset, distributed in 17 different clonal genotypes, with G21 being the most abundant in the population (376/910 ≈ 41%). Also in this case, samples collected in the initial weeks belong to the ancestral clades, supporting the hypothesis of multiple introductions. Notably, after the 17th week, all American samples display the haplotype g.23403A>G (S,p.614D>G) and we notice an overall decrease in genomic diversity, since the observed clonal genotypes pass from 16 (week interval 9–16, 2020) to 8 (week interval 17–29, 2020). Notice that only 49.1% of the American samples have a collection date.

Two-hundred and sixty samples from South Africa are included in the dataset, which are partitioned in six different clonal genotypes, four of which (G1, G8, G14, and G16) include

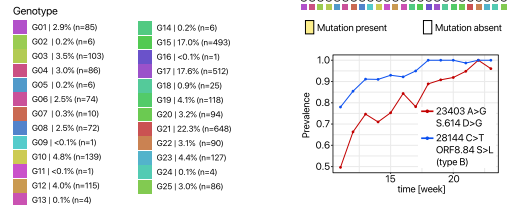
A VERSO | Viral Evolution ReconStruction



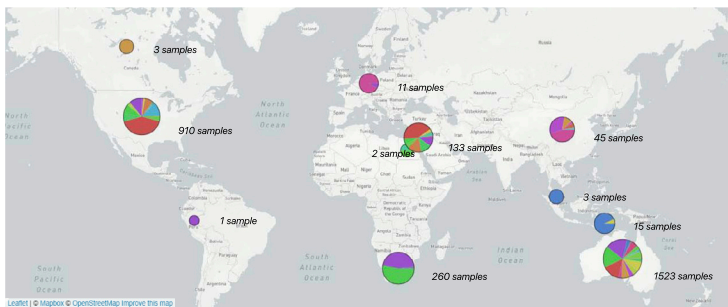
B

DATASET #1
Amplicon Illumina sequencing
n = 2906 samples
m = 29 clonal variants used in STEP #1
10571 variants detected in at least 1 sample

Variant	Genome location	Variant effect	Amino Acid substitution	Corrected clonal genotypes
29095 T>C	N	S	-	
18060 C>T	ORF1ab	S	-	
17858 A>G	ORF1ab	NS	S865 Y>C	
9477 T>A	ORF1ab	NS	3071 F>Y	
28657 C>T	N	S	-	
28144 C>T	ORF8	NS	84 S>L	
8782 T>C	ORF1ab	S	-	
11083 G>T	ORF1ab	NS	3606 L>F	
28311 C>T	N	NS	13 P>L	
28144 G>T	ORF3a	NS	291 G>V	
14505 C>T	ORF1ab	S	-	
28688 T>C	N	S	-	
29742 G>T	3'UTR	NC	-	
1397 G>A	ORF1ab	NS	378 V>I	
3037 C>T	ORF1ab	S	-	
241 C>T	5'UTR	NC	-	
14408 C>T	ORF1ab	NS	4715 P>L	
23403 A>G	S	NS	614 D>G	
28881 G>A	N	NS	203 R>K	
28882 G>A	N	S	-	
28883 G>C	N	NS	204 G>R	
25563 G>T	ORF3a	NS	57 Q>H	
2416 C>T	ORF1ab	S	-	
18877 C>T	ORF1ab	S	-	
1059 C>T	ORF1ab	NS	265 T>I	
11916 C>T	ORF1ab	NS	3884 S>L	
27964 C>T	ORF8	NS	24 S>L	
25572 C>T	ORF3a	S	-	
18555 C>T	ORF1ab	S	-	



C



D

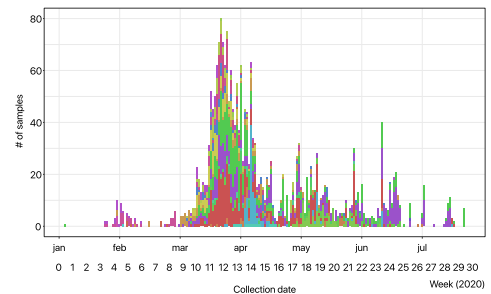


Figure 3. Viral evolution and intra-host genomic characterization of 2906 SARS-CoV-2 samples of via VERSO (dataset #1)

(A) The phylogenetic model returned by VERSO step #1 from the mutational profile of 2,906 samples selected after the quality check, on 29 clonal variants (VF > 90%) detected in at least 3% of the samples of dataset #1 (reference genome: SARS-CoV-2-ANC). Colors mark the 25 distinct clonal genotypes identified by VERSO (the mapping with the lineage nomenclature proposed in Rambaut et al.⁷⁸ and generated via pangolin 2.0⁷⁹ is provided in File S3). Samples with identical corrected clonal genotypes are grouped in polytomies and the black sample represents the SARS-CoV-2-ANC genome (visualization via FigTree⁸⁰). The green curves juxtaposed to certain polytomies report the number and fraction of samples in which the five homoplastic mutations are observed (only if the mutation is detected in at least 10 samples with the same corrected clonal genotype; see Data S2 for a summary on the samples exhibiting homoplastic clonal variants). The projection of the intra-host genomic diversity computed by VERSO step #2 from VF profiles is shown on the UMAP low-dimensional space for the clonal genotypes including ≥ 100 samples. Samples are clustered via Leiden algorithm on the kNN graph ($k = 10$), computed on the Bray-Curtis dissimilarity on VF profiles, after PCA. Solid lines represent the edges of the k-NNG.

(B) The composition of the corrected clonal genotypes returned by VERSO step #1 is shown. Clonal SNVs are annotated with mapping on ORFs, synonymous (S), nonsynonymous (NS), and non-coding (NC) states, and related amino acid substitutions. Variants g.8782T>C (*ORF1ab*, synonymous) and g.28144C>T (*ORF8*, p.84S>L) are colored in blue, whereas variant g.23403 A_IG (S, p.614 D_IG) is colored in red. The prevalence variation in time of the relative haplotypes (i.e., the fraction of samples displaying such mutations) is also shown. The five homoplastic variants are colored in green.

(C and D) (C) The geo-temporal localization of the clonal genotypes via Microreact⁸¹ and (D) the prevalence variation in time are displayed.

a single sample, whereas 98.46% of the samples exhibit the haplotype g.23403A>G (S,p.614D>G) and, specifically, are included in clonal genotypes G15 and G17. Finally, all Chinese samples were collected in the early phase (January–February, 2020) and are characterized by six different clonal genotypes (i.e., G1, G6, G7, G9, G11, and G12).

Homoplasmy detection (clonal variants)

Five clonal variants included in our model show apparent violations of the accumulation hypothesis, namely g.11083G>T (*ORF1ab*, p.3606 L>F), g.14805C>T (*ORF1ab*, synonymous), g.18555C>T (*ORF1ab*, synonymous), g.27964C>T (*ORF8*, p.24S>L), and g.28311C>T (N,p.13P>L), suggesting that they

might be involved in homoplasies. In [Figure 3](#) the samples in which the five homoplastic variants are detected are highlighted (if the mutation is detected in ≥ 10 samples with the same corrected clonal genotype), whereas in [Figure S3](#) one can find the expanded clonal variant tree, in which the reticulation related to such variants is explicitly depicted.

Some of such variants have been exhaustively studied (e.g., g.11083G>T in van Dorp et al.⁸⁹), specifically to verify possible scenarios of convergent evolution, which may unveil the fingerprint of adaptation of SARS-CoV-2 to human hosts. To this end, particular attention should be devoted to the three non-synonymous substitutions; i.e., g.11083G>T (present in 460 samples, $\approx 16\%$ of the dataset), g.27964C>T (182 samples, $\approx 6\%$) and g.28311C>T (153 samples, $\approx 5\%$). As a first result, we note the prevalence dynamics of the haplotypes defined by such variants does not show any apparent growth trend in the population (see [Figure S5](#)).

To further investigate if such variants fall in a region prone to mutations of the SARS-CoV-2 genome, we evaluated the mutational density employing a sliding window approach similarly to Soares et al.⁹⁰ (see [Supplemental experimental procedures](#) for additional details). As shown in [Figure S4](#), the mutational density, computed by considering synonymous minor variants, exhibits a median value of $= 0.083 [\text{syn.mutations}][\text{nucleotides}]^{-1}$. Interestingly, the three nonsynonymous SNVs (g.11083G>T, g.27964C>T and g.28311C>T) are located within windows with a higher mutational density than the median value: 0.085, 0.124, and 0.1 $\frac{\text{syn.mutations}}{\text{nucleotides}}$, respectively (see [Table S5](#)), and this would suggest that they might have originally emerged due to the presence of natural mutational hotspots or phantom mutations.

However, this analysis is not conclusive and further investigations are needed to characterize the functional effect of such mutations and the possible impact in the evolutionary and diffusion process of SARS-CoV-2.

Stability analysis

The choice of an appropriate VF threshold to identify clonal variants and, accordingly, to generate consensus sequences from raw sequencing data might affect the stability of the results of any downstream phylogenomic analysis. On the one hand, loose thresholds might increase the risk of including non-clonal variants in consensus sequences. On the other hand, too strict thresholds might increase the rate of false-negatives, especially with noisy sequencing data.

For this reason, we assessed the robustness of the results produced by VERSO step #1 on dataset #1 when different thresholds in the set $\delta \in \{0.5, 0.6, 0.7, 0.8\}$ are employed to identify clonal variants, with those obtained with default threshold ($\delta = 0.9$), in terms of tree accuracy (see the [Supplemental experimental procedures](#) for further details). As one can see in [Figure S7](#), the tree accuracy varies between 0.97 and 0.98 in all settings, proved the results produced by VERSO step #1 are robust with regard to the choice of the VF threshold for clonal variant identification.

VERSO step #2: Characterization of intra-host genomic diversity

We then applied VERSO step #2 to the complete VF profiles of the samples with the same clonal genotype and projected their intra-host genomic diversity on the UMAP low-dimensional space. This was done excluding (1) the clonal variants employed in the phylogenetic inference via VERSO step #1, (2) all minor var-

iants ($VF \leq 90\%$) observed in more than one clonal genotype (i.e., homoplasies) and that are likely emerged independently within the hosts, due to mutational hotspots, phantom mutations, or positive selection (see [Experimental procedures](#) and the next subsections). Even though, as expected, the VF profiles of minor variants are noisy, a complex intra-host genomic architecture is observed in several individuals. Moreover, patterns of co-occurrence of minor variants across samples support the hypothesis of transmission from one host to another.

In [Figure 3](#) we display the UMAP plots for the clonal genotypes including more than 100 samples, plus clonal genotype G4 ($n = 86$ samples), which was used for contact tracing analyses. Such maps describe likely transmission paths among hosts characterized by the same (corrected) clonal genotype and, in most cases, suggests the existence of several distinct infection clusters with different size and density. This result was achieved by exploiting the different properties of clonal and minor variants via the two-step procedure of VERSO.

Contact tracing

To corroborate our findings, we employed the contact tracing data from Rockett et al.,⁶⁹ in which 65 samples from dataset #1 are characterized with respect to household, work location, or other direct contacts. Four distinct contact groups, including 36, 15, 12, and 2 samples, respectively, are associated directly or indirectly to three different New South Wales institutions (i.e., institutions #1, #2, and #3) and to the same household environment (household #1).

As a first result, all samples belonging to a specific contact group are characterized by the same clonal genotype, determined via VERSO step #1, a result that confirms recent findings.^{42,69} More importantly, the analysis of the intra-host genomic diversity via VERSO step #2 allows one to highly refine this analysis.

In [Figure 4](#) one can find the UMAP plot of clonal genotypes G4, G12, and G21, which include 36 (over 86), 14 (over 115), and 14 (over 648) samples with contact information. Strikingly, the distribution of the pairwise intra-host genomic distance among samples from the same institution/household (computed on the K-nearest neighbor graph [k-NNG] via Bray-Curtis dissimilarity, after principal component analysis [PCA]; see [Experimental procedures](#)) is significantly lower with respect to the distance of all samples with the same clonal genotype (p value of the Mann-Whitney U test < 0.001 in all cases). Furthermore, all samples belonging to the same contact group are connected in the k-NNG, while a noteworthy proportion of samples without contact information in genotypes G12 and G21 are placed in disconnected graphs (24.9% and 76.4%, respectively).

This major result suggests that patterns of co-occurrence of minor variants can indeed provide useful indication on contact tracing dynamics, which would be masked when employing consensus sequencing data. Accordingly, the algorithmic strategy employed by VERSO step #2 and, especially, the identification of the k-NNG on intra-host genomic similarity provides an effective tool to dissect the complexity of viral evolution and transmission, which might in turn improve the reliability of currently available contact tracing tools.

Homoplasy detection (minor variants)

Several minor variants are found in samples with distinct clonal genotypes and might indicate the presence of homoplasies. In

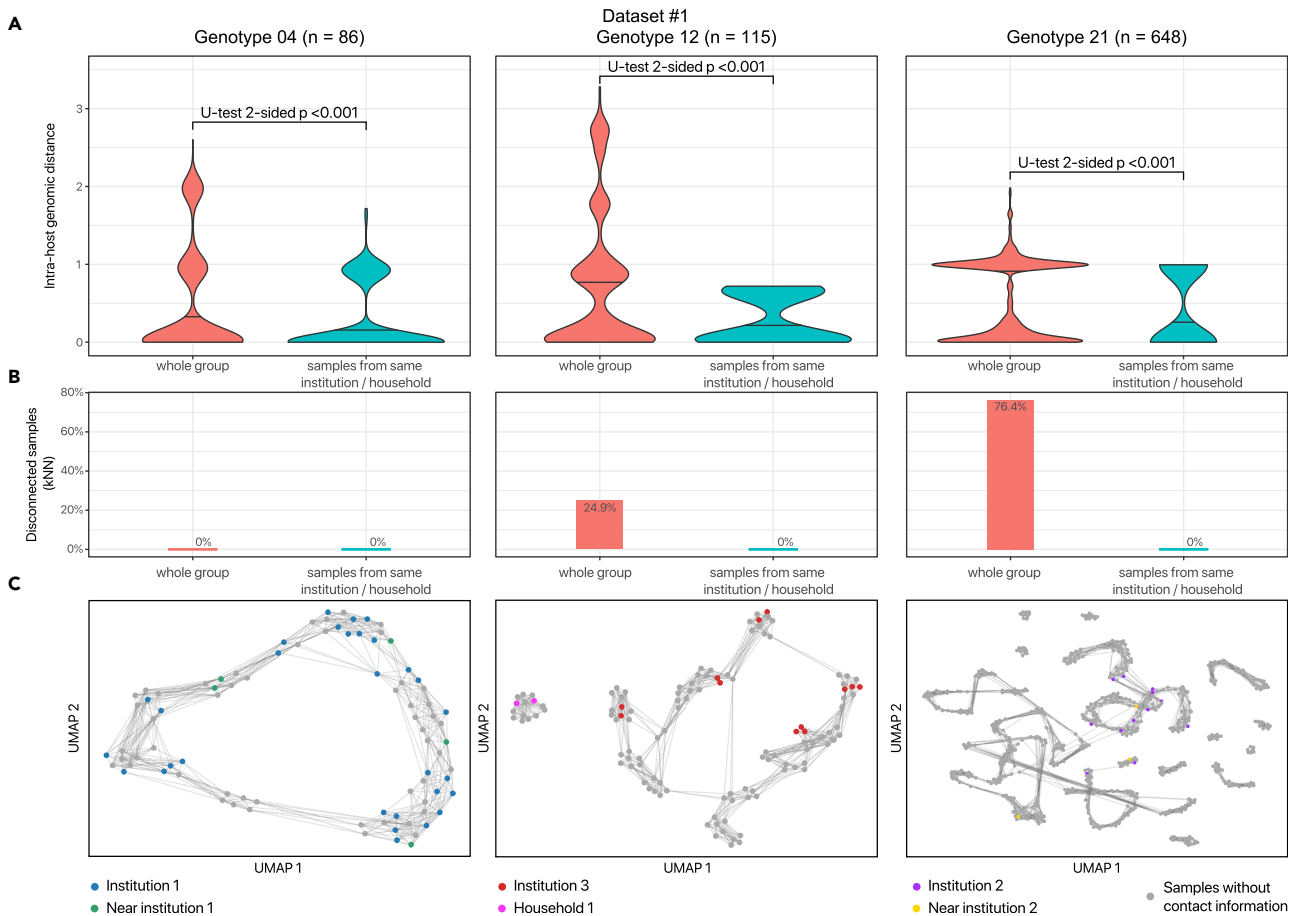


Figure 4. Infection dynamics revealed via characterization of intra-host genomic similarity (dataset #1)

(A) The distribution of the pairwise intra-host genomic distance (computed via Bray-Curtis dissimilarity on the kNN graph, with $k = 10$, after PCA; see [Experimental procedures](#)) for the samples belonging to the same household or institution (including samples marked as near), versus the pairwise distance of all samples belonging to clonal genotypes G4, G12, and G21. The p values of the Mann-Whitney U test two-sided are also shown.

(B) The proportion of samples that are disconnected in the kNN graph, with respect to the samples belonging to the same household or institution (including samples marked as near) and with respect to all samples.

(C) The UMAP projection of the intra-host genomic diversity of the samples belonging to clonal genotypes G4, G12, and G21, returned by VERSO step #2.

this respect, the heatmap in [Figure 5F](#) returns the distribution of minor SNVs with respect to (1) the number of distinct clonal genotypes in which they are detected, and (2) the mutational density of the region in which they are located (see the [Supplemental experimental procedures](#) for details on the mutational density analysis).

The intuition is that the variants detected in single clonal genotypes (left region of the heatmap) are likely spontaneously emerged private mutations, or the result of infection events between hosts with same clonal genotype (see above). Conversely, SNVs found in multiple clonal genotypes (right region of the heatmap) may have emerged due to positive selection in a parallel/convergent evolution scenario, or to mutational hotspots or phantom mutations. To this end, the mutational density analysis provides useful information to pinpoint mutation-prone regions of the genome.

Interestingly, a significant number of minor variants are observed in multiple clonal genotypes and fall in scarcely mutated regions of the genome (see [Figure S4](#)). This would suggest that some of these variants might have been positively

selected, due to some possible functional advantage or to transmission-related founder effects. In this respect, we further focused our investigation on a list of 80 candidate minor variants that (1) are detected in more than one clonal genotype, (2) are present in at least 10 samples, (3) are nonsynonymous, and (4) fall in a region of the genome with mutational density lower than the median value (see [Table S6](#) for details on such variants). In the following, we focus on a subset of such variants falling on the spike gene of the SARS-CoV-2 genome.

Considerations on homoplasies falling on the spike gene

The spike protein of SARS-CoV-2 plays a critical role in the recognition of the ACE2 receptor and in the ensuing cell membrane fusion process.⁹¹ We prioritized three candidate homoplasmic minor variants occurring on the SARS-CoV-2 spike gene (S) (see [Table S6](#)). Interestingly, two out of three, namely $g.24552T>C$ ($p.997I>T$) and $g.24557G>T$ ($p.999G>C$), detected in 57 samples in total (10 and 47 samples, respectively), clustered in the so-called connector region (CR), bridging between the two heptad repeat regions (HR1 and HR2) of the S2 subunit of the spike protein.

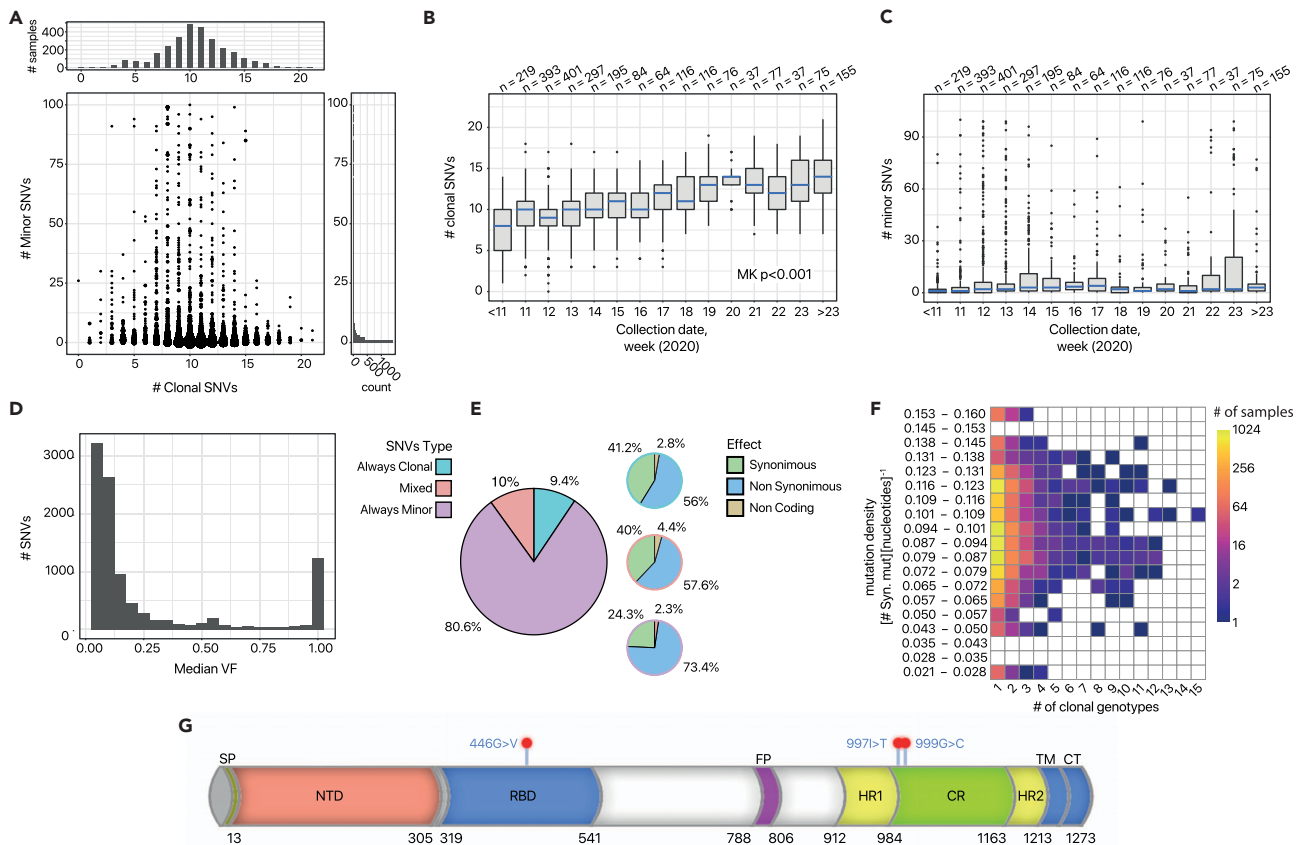


Figure 5. Mutational landscape of 2906 SARS-CoV-2 samples (dataset #1)

(A) Scatterplot displaying, for each sample, the number of clonal ($VF > 90\%$) and minor variants ($VF \leq 90\%$, node size proportional to the number of samples). (B and C) Boxplots returning the distribution of the number of clonal (B) and minor variants (C), obtained by grouping samples according to collection date (weeks, 2020). The p value of the Mann-Kendall (MK) trend test on clonal variants is highly significant. (D) Distribution of the median VF for all SNVs detected in the viral populations. (E) Pie charts returning (left) the proportion of SNVs detected as always clonal, always minor, or mixed; (right) for each category, the proportion of synonymous, nonsynonymous, and non-coding variants (check the pie-chart border color for a visual clue). (F) Heatmap returning the distribution of always minor SNVs with respect to (x axis) the number of clonal genotype of the phylogenomic model in Figure 3 in which each variant is observed, (y axis) the mutational density of the genome region in which it is located (see the Supplemental experimental procedures). (G) Mapping of the candidate homoplastic minor variants located on the spike gene of the SARS-CoV-2 virus.

When the receptor binding domain (RBD) binds to ACE2 receptor on the target cell, it causes a conformational change responsible for the insertion of the fusion peptide (FP) into the target cell membrane. This, in turn, triggers further conformational changes, eventually promoting a direct interaction between HR1 trimer and HR2, which occurs upon bending of the flexible CR, in order to form a six-helical HR1-HR2 complex known as the fusion core region (FCR) in close proximity to the target cell plasma membrane, ultimately leading to viral fusion and cell entry.⁹²

Peptides derived from the HR2 heptad region of enveloped viruses and able to efficiently bind to the viral HR1 region inhibit the formation of the FCR and completely suppress viral infection.⁹³ Therefore, the formation of the FCR is considered to be vital to mediate virus entry in the target cells, promoting viral infectivity. Of note, the CR is highly conserved across the Gammacoronavirus genus, supporting the notion that this region may play a very important but still unclear functional role (Figure 5G). Although structural and *in vitro* models will be required in order to exten-

sively characterize the functional effect of these variants, the evidence that two of our three minor variants detected in the spike protein falls in a small domain comprising less than 14% of the entire spike protein length is intriguing, as it suggests a potential functional role for these mutations. It will be important to track the prevalence of these mutations, as well as of all other candidate convergent variants falling on different region of the SARS-CoV-2, to highlight possible transitions to clonality (see below). We also remark that, being a data-science computational approach, VERSO can struggle in dissecting complex mutational cases, since all the experimental hypotheses that can be generated are clearly data dependent. For this reason, and given the heterogeneity and limitations of currently available SARS-CoV-2 datasets, any hypothesis delivered by VERSO requires additional independent investigations and *ad hoc* experimental validations.

Mutational landscape

We analyzed in depth the mutational landscape of the samples of dataset #1. First, the comparison of the number of clonal

(VF > 90%) and minor variants detected in each host (Figure 5A) reveals a bimodal distribution of clonal variants (with first mode at 4 and s mode at 10), whereas minor variants display a more dispersed long-tailed distribution with median equal to 2 and average ≈ 23 . From the plot, it is also clear that individuals characterized by the same clonal genotype may display a significantly different number of minor variants, with distinct distributions observed across clonal genotypes.

The comparison of the distribution of the number of variants obtained by grouping the samples with respect to collection week (Figures 5B and 5C) allows us to highlight a highly statistically significant increasing trend for clonal variants (Mann-Kendall trend test on median number of clonal variants, $p < 0.001$). This result would strongly support both the hypothesis of accumulation of clonal variants in the population and that of a concurrent increase of overall genomic diversity of SARS-CoV-2,^{36,94} whereas the relevance of this phenomenon on minor variants is unclear.

We then focused on the properties of the SNVs detected in the population. Surprisingly, the distribution of the median VF for each detected variant (Figure 5D) reveals a bimodal distribution, with the large majority of variants showing either a very low or a very high VF, with only a small proportion of variants showing a median VF within the range 10%–90%. This behavior is typical of systems where the prevalence of some subpopulations is driven by positive Darwinian selection while others are purified.⁹⁵

In order to analyze the two components of this distribution, we further categorized the variants as always clonal (i.e., SNVs detected with VF >90% in all samples), always minor (i.e., SNVs detected with VF 5% and $\leq 90\%$ in all samples), and mixed (i.e., SNVs detected as clonal in at least one sample and as minor in at least another sample). As one can see in Figure 5E, 9.4%, 80.6%, and 10% all SNVs are respectively detected as always clonal, always minor, and mixed in our dataset. Moreover, 56%, 73.4%, and 57.6% of always clonal, always minor, and mixed variants, respectively, are nonsynonymous, whereas the large majority of remaining variants are synonymous.

These results would suggest that, in most cases, randomly emerging SARS-CoV-2 minor variants tend to remain at a low frequency in the population, whereas, in some circumstances, certain variants can undergo frequency increases and even become clonal, due to undetected mixed transmission events or to selection shifts, as it was observed by Poon et al.⁸ for the cases of H3N2 and H1N1/2009 influenza. Interestingly, 15 variants identified as possibly convergent (see above) fall into this category and deserve further investigations (see Table S6 for additional details).

Transmission bottleneck analysis

The estimation of transmission bottlenecks might be of specific interest during the current pandemics. Despite most available methods requiring data collected on donor-host couples (see, e.g., Sobel Leonard et al.⁹⁶ and Ghafari et al.⁹⁷), here we employed a strategy akin to Monsion et al.⁹⁸ and Lequime et al.⁹⁹ that is roughly based on the analysis of the variation of the VF variance of a number of candidate neutral mutations. The intuition is that variance shrinking indicates significant transmission bottlenecks, which, accordingly, would result in lower viral diversity transferred from a host to another and, possibly, in purification of certain variants in the population. As the analysis ideally

requires the comparison of groups in which infection events have occurred, here we considered groups of samples with distinct clonal genotypes, separately. We then selected a number of variants as neutral markers. The rationale is that transmission phenomena such as bottlenecks are expected to significantly affect the VF variance of neutral markers (please see Supplemental experimental procedures for further details).

More in detail, we first split the samples of each clonal genotype for which a collection date is available into non-overlapping groups corresponding to two consecutive time windows; i.e., before and after the 14th week, 2020. Accordingly, three SNVs were selected as candidate-neutral or quasineutral markers, namely variants g.634T>C, g.14523A>G, and g.15168G>A. In Figure S6, one can find the distribution of the VF of the selected markers with respect to the time windows, which highlights moderate variations of the variance for all markers (see also Table S7. All in all, this result would suggest the presence of mild bottleneck effects, consistent with recent studies involving donor-host data.⁴³

Application of VERSO to 2,766 samples from RNA-sequencing data (dataset #2)

We retrieved the raw Illumina RNA-sequencing data of 2,766 samples included in dataset #2 and applied VERSO to the mutational profiles of 1,438 samples selected after quality check. Twenty-three clonal variants were employed in the analysis, according to the filters described later.

The resulting phylogenetic model is consistent with the one obtained for dataset #1, despite minor differences (Figure S8A). Specifically, 18 distinct clonal genotypes are identified by VERSO step #1, 11 of which are identical to those found in the analysis of dataset #1 (in such cases the same genotype label was maintained; see Data S3 for the mapping with the lineage nomenclature proposed by Rambaut et al.⁷⁸). Five further clonal genotypes are evolutionarily consistent and represent independent branches detected due to the non-overlapping composition of the dataset, and are labeled with progressive letters from the closest genotype (i.e., G13a, G21a, G22a, G22b, G23a), while the two samples of genotype G13a* might be safely assigned to genotype G13a, since the absence of mutation g.3037C>T is likely due to low coverage.

By excluding the remaining clonal genotype GH, which presents inconsistencies due to the presence of the candidate homoplastic variant g.11083G>T (*ORF1ab*, p.3606L>F, see above), all clonal genotypes display the same ordering in both datasets (also see the expanded clonal variant tree in Figure S9). This proves the robustness of the results delivered by VERSO step #1 even when dealing with data generated from distinct sequencing platforms.

By looking at the geo-temporal localization of samples obtained via Microreact⁸¹ (Figure S8B), one can see that that dataset #2 includes samples with a significantly different geographical distribution with respect to dataset #1. This dataset contains sample from 10 countries, with the large majority collected in the United States (96.8%). More in detail, the samples of such countries are mostly characterized by clonal genotype G21. We further notice that, also for dataset #2, mutation g.23403A>G (S,p.614D>G) becomes prevalent in the population at late collection dates. Moreover, only samples belonging to previously defined type B are detected in this dataset.

The analysis of the intra-host genomic diversity was also performed for dataset #2 via VERSO step #2, which would suggest the existence of undetected infection events and of several infection clusters with distinct properties, even though no contact tracing is available in this case. Overall, this proves the general applicability of the VERSO framework, which can produce meaningful results when applied to data produced with any sequencing platforms. However, in order to minimize the possible impact of data- and platform-specific biases, our suggestion is to perform the VERSO analysis on datasets generated from different protocols separately.

Scalability

We finally assessed the computational time required by VERSO in a variety of simulated scenarios. The results are shown in the [Supplemental experimental procedures \(Figure S10\)](#) and demonstrate the scalability of VERSO also when processing large-scale datasets.

DISCUSSION

We introduced VERSO, a comprehensive framework for the high-resolution characterization of viral evolution from sequencing data, which is an improvement on currently available methods for the analysis of consensus sequences. VERSO exploits the distinct properties of clonal and minor variants to dissect the complex interplay of genomic evolution within hosts and transmission among hosts.

On the one hand, the probabilistic framework underlying VERSO step #1 delivers highly accurate and robust phylogenetic models from clonal variants, also in conditions of noisy observations and sampling limitations, as proved by extensive simulations and by the application to two large-scale SARS-CoV-2 datasets generated from distinct sequencing platforms. On the other hand, the characterization of intra-host genomic diversity provided by VERSO step #2 allows one to identify undetected infection paths, which were in our case validated with contact tracing data, as well as to intercept variants involved in homoplasies.

This may represent a major advancement in the analysis of viral evolution and spread and should be quickly implemented in combination with data-driven epidemiological models to deliver a high-precision platform for pathogen detection and surveillance.^{12,100} This might be particularly relevant for countries that suffered outbreaks of exceptional proportions and for which the limitations and inhomogeneity of diagnostic tests have proved insufficient to define reliable descriptive/predictive models of disease diffusion. For instance, it was hypothesized that the rapid diffusion of COVID-19 might be likely due to the extremely high number of untested asymptomatic hosts.¹⁰¹

More accurate and robust phylogenetic models may allow one to improve the assessment of molecular clocks and, accordingly, the estimation of the parameters of epidemiological models such as susceptible-infected-recovered (SIR) and susceptible-infected-susceptible (SIS),^{11,102} as well as to unravel the cryptic transmission paths.^{8,12,13,103} Furthermore, the finer grain of the analysis on intra-host genomic similarity from sequencing data might be employed to enhance the active surveillance; for instance, by facilitating the identification of infection clusters and super-spreaders.¹⁰⁴ Finally, the

characterization of variants possibly involved in positive selection processes might be used to drive the experimental research on treatments and vaccines.

EXPERIMENTAL PROCEDURES

Resource availability

Lead contact

Alex Graudenzi, Institute of Molecular Bioimaging and Physiology, Consiglio Nazionale delle Ricerche (IBFM-CNR), via F.lli Cervi, 93, 20,090 Segrate, Milan, Italy. alex.graudenzi@ibfm.cnr.it.

Materials availability

This study did not generate new unique reagents.

Data and code availability

VERSO is freely available at this link: <https://github.com/BIMIB-DISCO/VERSO>. VERSO step #1 is provided as an open source standalone R tool, whereas step #2 is provided as a Python script. The source code to replicate all the analyses presented in the manuscript, both on simulated and real-world datasets, is available at this link: <https://github.com/BIMIB-DISCO/VERSO-UTILITIES>.

SCANPY⁵⁹ is available at this link: <https://scanpy.readthedocs.io/en/stable/>. The Web-based tool for the geo-temporal visualization of samples, Microreact,⁸¹ is available at this link: <https://microreact.org/showcase>. The tool employed to plot the phylogenomic model returned by VERSO step #1 (in Newick file format) is FigTree⁸⁰ and is available at this link: <http://tree.bio.ed.ac.uk/software/figtree/>. The tool used for the mapping between clonal genotype labels and the dynamic nomenclature proposed by Rambaut et al.⁷⁸ is pangolin 2.0⁷⁹ and is available at this link: <https://github.com/cov-lineages/pangolin>.

VERSO step #1: robust phylogenomic inference from clonal variant profiles

VERSO is a novel framework for the reconstruction of viral evolution models from raw sequencing data of viral genomes. It includes a two-step procedure, which we describe in the following.

The first step of VERSO employs a probabilistic maximum-likelihood framework for the reconstruction of robust phylogenetic trees from binarized mutational profiles of clonal variants (or, alternatively, from consensus sequences). This step relies on an evolved version of the algorithmic framework introduced by Ramazzotti et al.¹⁰⁵ for the inference of cancer evolution models from single-cell sequencing data, and can be executed independently from step #2, in case raw sequencing data are not available.

Inputs

The method takes as input a n (samples) \times m (variants) binary mutational profile matrix, as defined on the basis of clonal SNVs only. In this case, an entry in a given sample is equal to 1 (present) if the VF is larger than a certain threshold (in our analyses, equal to 90%), it is equal to 0 if lower than a distinct threshold (in our analyses, equal to 5%), and is considered as missing (NA) in the other cases, thus modeling possible uncertainty in sequencing data or low coverage.

Notice that consensus sequences can be processed by VERSO step #1 by generating a consistent binarized mutational profile matrix. We also note that, given the intrinsic challenges associated with a reliable identification of low VF indels, the analysis focuses only on SNVs. Further details on the variant calling pipeline employed in this study are provided next.

The algorithmic framework

VERSO step #1 is a probabilistic framework that solves a Boolean matrix factorization problem with perfect phylogeny constraints and relying on the infinite sites assumption (ISA).^{106,107} The ISA subsumes a consistent process of accumulation of clonal variants characterizing the evolutionary history of the virus and does not allow for losses of mutations or convergent variants (i.e., mutations observed in distinct clades).

In this regard, we recall that the variant accumulation hypothesis holds only when considering clonal mutations. In fact, clonal mutations (e.g., A, B, C, D, and E) are present, by definition, in the large majority of the quasispecies of a given sample, depending on the chosen VF threshold (in our case, equal to 90%; see above). Since such variants are rarely lost, they are most likely transmitted from one host to another during infections. In addition, the origination of

new clonal mutations in single samples leads to the definition of new clonal genotypes, following a standard branching process (e.g., A, AB, ABC, ABD, ABDE). As a result, clonal mutations typically accumulate during the evolutionary history of a virus, excluding complex scenarios involving reticulation events,⁵³ whereas clonal genotypes can clearly become extinct. Conversely, variants with lower frequency do not necessarily accumulate, due to the high recombination rates, as well as to bottlenecks, founder effects, and stochasticity,³¹ and this is the reason why they were considered separately in the analysis, via VERSO step #2 (see below).

More in detail, VERSO step #1 accounts for uncertainty in the data, by employing a maximum-likelihood approach (via Markov chain Monte Carlo [MCMC] search) that allows for the presence of false-positives, false-negatives, and missing data points in clonal variant profiles. As shown by Ramazzotti et al.¹⁰⁵ in a different experimental context, our algorithmic framework ensures robustness and scalability also in case of high rates of errors and missing data, due, for instance, to sampling limitations. Furthermore, it is robust to mild violations of the ISA (e.g., due to reticulation events, such as convergent variants) or mutation losses, which can be characterized after the inference, if present (see the specific features on homoplasy detection discussed next). Please refer to the [Supplemental experimental procedures](#) for further details on the algorithmic framework and its assumptions, including the probabilistic graphical model depicted in [Figure S1](#) and the summary of notation in [Table S4](#).

Outputs

The inference returns a set of maximum-likelihood variants trees (minimum 1) as sampled during the MCMC search, representing the ordering of accumulation of clonal variants, and a set of maximum-likelihood attachments of samples to variants. Given the variants tree and the maximum-likelihood attachments of samples to variants, VERSO outputs (1) a phylogenetic model where each leaf correspond to a sample, whereas internal nodes correspond to accumulating clonal variants; (2) the corrected clonal genotype of each sample (i.e., the binary mutational profile on clonal variants obtained after removing false-positives, false-negatives, and missing data).

The model naturally includes polytomies, which group samples with the same corrected clonal genotype. The length of the branches in the model represents the number of clonal substitutions (which can be normalized with respect to genome length), as in standard phylogenomic models, and the clades of the model correspond to viral lineages. The VERSO phylogenetic model is provided as output in Newick file format and can be processed and visualized in standard tools for phylogenetic analysis, such as FigTree⁸⁰ or Dendroscope.¹⁰⁸ Furthermore, VERSO allows one to visualize the geo-temporal localization of clonal genotypes via Microreact.⁸¹

Additional feature: homoplasy detection on clonal variants

Violations of the ISA are possible and can be due to reticulation events⁵³ such as homoplasies (i.e., identical variants detected in samples belonging to different clades) or to rare occurrences involving mutation losses (e.g., due to recombination-related deletions or to multiple mutations hitting an already mutated genome location³⁴), as well as to infrequent transmission phenomena, such as super-infections^{65,66} (a discussion on the general limitations of approaches based on phylogenetic trees when dealing with reticulation events is available elsewhere^{109–111}).

In this regard, VERSO allows one to identify clonal mutations likely involved in homoplasies, in a similar fashion to the plethora of works on mitochondrial evolution and phylogenetic networks (discussed elsewhere^{53,54,56,112–114}). In detail, given the maximum-likelihood phylogenetic tree, VERSO can estimate the variants that are theoretically expected in each sample. By comparing the theoretical observations with the input data, VERSO can estimate the rate of false-positives (i.e., the variants that are observed in the data but are not predicted by VERSO), and false-negatives (i.e., variants that are not observed but predicted). Variants that show particularly high estimated error rates represent candidate homoplasies and are flagged. First, this allows one to pinpoint samples exhibiting homoplastic mutations (see [Figures 3](#) and [S8](#)) and, second, to reconstruct an expanded clonal variants tree, in which candidate homoplastic mutations are duplicated after the inference, so to allow the visualization of reticulation events, as proposed by Skála and Zrzavý¹¹² (see, e.g., [Figures S3](#) and [S9](#)).

Furthermore, once this procedure has been completed, the list of flagged variants can include (1) mutations falling in highly mutated regions due to muta-

tional hotspots, (2) phantom mutations (i.e., systematic artifacts generated during sequencing processes⁵⁶), or (3) mutations that have been positively selected in the population (e.g., due to a particular functional advantage).

Since one might be interested in identifying positively selected mutations, VERSO allows one to perform a consecutive analysis that aims at highlighting the mutation-prone regions of the genome and that might be due to mutational hotspots or phantom mutations (see the [Supplemental experimental procedures](#) for further details). We finally note that the detection of homoplasies for minor variants requires a different algorithmic procedure, which is detailed in the following.

VERSO step #2: characterization of intra-host genomic diversity

In the second step, VERSO takes into account the VF profiles of groups of samples with the same corrected clonal genotype (identified via VERSO step #1), in order to characterize their intra-host genomic diversity and visualize it on a low-dimensional space. This allows one to highlight patterns of co-occurrence of minor variants, possibly underlying undetected infection events, as well as homoplasies involving; e.g., positively selected variants. Notice that this step requires raw sequencing data and the prior execution of step #1.

Inputs

VERSO step #2 takes as input a n (samples) \times m (variants) VF profile matrix, in which each entry includes the $VF \in (0, 1)$ of a given mutation in a certain sample, after filtering out (1) the clonal variants employed in step #1 and (2) the minor variants possibly involved in homoplasies (see below). The variant calling pipeline employed in this work is detailed next.

The algorithmic framework

While it is sound to binarize clonal variant profiles to reconstruct a phylogenetic tree, it is opportune to consider the VF profiles when analyzing intra-host variants, for several reasons. First, VF profiles describe the intra-host genomic diversity of any given host, and this information would be lost during binarization. Second, minor variant profiles might be noisy, due to the relatively low abundance and to the technical limitations of sequencing experiments. Accordingly, such data may possibly include artifacts, which can be partially mitigated during the quality-check phase and by including in the analysis only highly confident variants. However, binarization with arbitrary thresholds might increase the false-positive rate, compromising the accuracy of any downstream analysis. Third, as specified above, the extent of transmission of minor variants among individuals is still partially obscure. The VF of minor variants is, in fact, highly affected by recombination processes, as well as by complex transmission phenomena, involving stochastic fluctuations, bottlenecks, and founder effects, which may lead certain variants changing their VF, not being transmitted, or even becoming clonal in the infected host.⁵⁷ The latter issue also suggests that the hypothesis of accumulation of minor variants during infections may not hold and should be relaxed.

For these reasons, VERSO step #2 defines a pairwise genomic distance, computed on the VF profiles, to be used in downstream analyses. The intuition is that samples displaying similar patterns of co-occurrence of minor variants might have a similar quasispecies architecture, thus being at a small evolutionary distance. Accordingly, this might indicate a direct or indirect infection event. In particular, in this work we employed the Bray-Curtis dissimilarity, which is defined as follows: given the ordered VF vectors of two samples (i.e. $v_i = \{VF_i^1, \dots, VF_i^m\}$ and $v_j = \{VF_j^1, \dots, VF_j^m\}$), the pairwise Bray-Curtis dissimilarity $d(i, j)$ is given by:

$$d(v_i, v_j) = \frac{\sum_{l=1}^m |VF_i^l - VF_j^l|}{\sum_{l=1}^m |VF_i^l + VF_j^l|} \quad (\text{Equation 1})$$

Since this measure weights the pairwise VF dissimilarity on each variant with respect to the sum of the VF of all variants detected in both samples, it can be effectively used to compare the intra-host genomic diversity of samples, as proposed, for instance, by Srinivas et al.¹¹⁵ However, VERSO allows one to employ different distance metrics on VF profiles, such as correlation or Euclidean distance.

As a design choice, in VERSO, the genomic distance is computed among all samples associated to any given corrected clonal genotype, as inferred in step #1. The rationale is that, in a statistical inference framework modeling a complex interplay involving heterogeneous dynamical processes, it is crucial to stratify samples into homogeneous groups, to reduce the impact of possible

confounding effects.¹¹⁶ Furthermore, as specified above, due to the distinct properties of clonal and minor variants during transmission, it is reasonable to assume that the event in which certain minor variants and no clonal variants are transmitted from a host to another during the infection is extremely unlikely. Accordingly, the clonal variants employed for the reconstruction of the phylogenetic tree in step #1 are excluded from the computation of the intra-host genomic distance among samples.

In order to produce useful knowledge from the genomic distance discussed above and since, in real-world scenarios, this is a typically complex high-dimensional problem, it is sound to employ state-of-the-art strategies for dimensionality reduction and (sample) clustering, as typically done in single-cell analyses.¹¹⁷ In this regard, the workflow employed in VERSO ensures high scalability with large datasets, also making it possible to take advantage of effective analysis and visualization features. In detail, the workflow includes three steps: (1) the computation of k-NNG, which can be executed on the original VF matrix, or after applying PCA, to possibly reduce the effect of noisy observations (when the number of samples and variants is sufficiently high); (2) the clustering of samples via either Louvain or Leiden algorithms for community detection;¹¹⁸ (3) the projection of samples on a low-dimensional space via standard tSNE⁶⁸ or UMAP⁶⁷ plots.

Outputs

As output, VERSO step #2 delivers both the partitioning of samples in homogeneous clusters and the visualization in a low-dimensional space, also allowing samples to be labeled according to other covariates, such as collection date or geographical location. In the map in Figure 3, for instance, the intra-host genomic diversity of each sample and the genomic distance among samples are projected on the first two UMAP components, whereas samples that are connected by k-NNG edges display similar patterns of co-occurrence of variants. Accordingly, the map shows clusters of samples likely affected by infection events, in which (a fraction of) quasiespecies might have been transmitted from one host to another. This represents a major novelty introduced by VERSO and also allows one to effectively visualize the space of VF profiles.

To facilitate the usage, VERSO step #2 is provided as a Python script which employs the SCANPY suite of tools,⁵⁹ which is typically used in single-cell analyses and includes a number of highly effective analysis and visualization features.

Additional feature: homoplasy detection on minor variants

Also in the case of minor variants, it is important to pinpoint possible homoplasies that might be due to mutational hotspots, phantom mutations, and convergent variants. Given the phylogenetic model retrieved via step #1, VERSO allows one to flag the variants that are detected in a number of clonal genotypes exceeding a user-defined threshold. In our case, the threshold is equal to 1, meaning that all minor variants found in more than one clonal genotype are flagged.

Such variants are then excluded from the computation of the intra-host genomic distance, prior to the execution of step #2. Furthermore, the list of flagged variants can be investigated as proposed for step #1 (see above), in order to possibly identify mutations involved in positive selection scenarios.

Datasets description

Dataset #1 (Illumina Amplicon sequencing)

We analyzed 3,960 samples from distinct individuals obtained from 22 NCBI BioProjects, which, at the time of writing, are all the publicly available datasets including raw Illumina Amplicon sequencing data. In detail, we selected the following NCBI BioProjects: (1) PRJNA613958, (2) PRJNA614546, (3) PRJNA616147, (4) PRJNA622817, (5) PRJNA623683, (6) PRJNA625551, (7) PRJNA627229, (8) PRJNA627662, (9) PRJNA629891, (10) PRJNA631042, (11) PRJNA633948, (12) PRJNA634119, (13) PRJNA636446, (14) PRJNA636748, (15) PRJNA639066, (16) PRJNA643575, (17) PRJNA645906, (18) PRJNA647448, (19) PRJNA647529, (20) PRJNA650037, (21) PRJNA656534, and (22) PRJNA656695.

Dataset #2 (Illumina RNA sequencing)

We analyzed 2,766 samples from distinct individuals obtained from 22 NCBI BioProjects, which, at the time of writing, are all the publicly available datasets including raw Illumina RNA-sequencing data. In detail, we selected the following NCBI BioProjects: (1) PRJNA601736, (2) PRJNA603194, (3) PRJNA605983, (4) PRJNA607948, (5) PRJNA608651, (6) PRJNA610428, (7) PRJNA615319, (8)

PRJNA616446, (9) PRJNA623895, (10) PRJNA624792, (11) PRJNA626526, (12) PRJNA631061, (13) PRJNA636446, (14) PRJNA637892, (15) PRJNA639591, (16) PRJNA639864, (17) PRJNA650134, (18) PRJNA650245, (19) PRJNA655577, (20) PRJNA657938, (21) PRJNA657985, and (22) PRJNA658211.

Contact tracing data

Contact tracing data were obtained from the study presented by Rockett et al.⁶⁹ In detail, for 65 samples included in dataset #1 (NCBI BioProject: PRJNA633948), information on households, work institutions, and epidemiological linkages are provided. Thus, it is possible to identify three different contact groups based on institutions regularly frequented by patients and one-household couples. Contact information was employed to assess the relation between the intra-host genomic similarity and the contact dynamics. The results are provided in the main text.

Parameter settings

Parameter settings of variant calling (datasets #1 and #2)

We converted all the samples to FASTQ files using the Sequence Read Archive (SRA) toolkit. Following Bastola et al.,⁷⁵ we used Trimmomatic (version 0.39) to remove the nucleotides with low quality score from the RNA sequences with the following settings: LEADING:20 TRAILING:20 SLIDINGWINDOW:4:20 MINLEN:40. We then used bwa mem (version 0.7.17) to map reads to SARS-CoV-2-ANC reference genome (Data S1; see Results). We generated sorted BAM files from bwa mem results using SAMtools (version 1.10) and removed duplicates with Picard (version 2.22.2). Variant calling was performed generating mpileup files using SAMtools and then running VarScan (min-var-freq parameter set to 0.01).¹¹⁹

We note that it was recently reported that some currently available SARS-CoV-2 datasets exhibit quality issues.^{13,120} Accordingly, one should be extremely careful when performing quality check and, especially, when considering low-frequency variants, which might possibly result from sequencing artifacts even in case of high-coverage experiments. In this regard, many effective approaches can be employed to reduce false variants. For instance, the Broad Institute recently updated an effective variant calling pipeline for viral genome data,¹²¹ while new methods for error correction of viral sequencing have been proposed at a widely used website (<https://virological.org>), which also includes a number of useful up-to-date guidelines and best practices for viral evolution analyses.

In our case, we here employed the following significance filters on variants. In particular, we kept only the mutations (1) showing a VarScan significance p value <0.01 (Fisher's exact test on the read counts supporting reference and variant alleles) and more than 25 reads of support in at least 75% of the samples, (2) displaying a VF >5%. As a result, we selected a list of 15,892 (over 55,280 overall SNVs) highly confident SNVs for dataset #1 and 7,389 (over 53,354) for dataset #2.

High-quality variants were then mapped on SARS-CoV-2 coding sequences (CDSs) via a custom R script, also by highlighting synonymous/nonsynonymous states and amino acid substitutions for the related open reading frame (ORF) product. In particular, we translated reference and mutated CDSs with the seqinr R package to obtain the relative amino acid sequences, which we compared to assess the effect of each nucleotide variation in terms of amino acid substitution.

We finally note that availability of the cycle threshold (Ct) values generated by qPCR and the related quantification of the amounts of viral transcripts would be very useful to characterize samples with high viral load, yet this information is not available for the considered datasets.

Quality check (datasets #1 and #2)

In order to select high-quality samples, we selected only those exhibiting high coverage and in particular those with at least 25 reads in more than 90% of the SARS-CoV-2-ANC genome. In addition, we filtered out all samples exhibiting more than 100 minor variants (VF ≤ 90%).

We finally excluded samples SRR11597146 and SRR11476447 from dataset #1, as the first sample displays zero SNVs and the second one reports an unfeasible collection date (i.e., 30th Jan. 2019).

After the quality-check filters, 2,906 samples of dataset #1 are left for downstream analyses, in which 10,571 distinct high-quality SNVs are observed, and 1,438 samples are left for dataset #2, with 6,143 high-quality SNVs.

Parameter settings of VERSO (datasets #1 and #2)

The phylogenomic analysis via VERSO step #1 was performed on datasets #1 and #2 by considering only clonal variants (VF > 90%) detected in at least 3% of the samples. A grid search comprising 16 different error rates was employed (see Table S3). Samples with the same corrected clonal genotype were grouped in polytomies in the final phylogenetic models.

The analysis of the intra-host genomic diversity via VERSO step #2 was performed by considering the VF profiles of all samples, by excluding (1) the clonal variants employed in the phylogenomic reconstruction via VERSO step #1, (2) the minor variants involved in homoplasies (i.e., observed in more than one clonal genotype returned by VERSO step #1). Missing values (NA) were imputed to 0 for downstream analysis. A number of principal components equals to 10 was employed in PCA step, prior to the computation of the k-NNG ($k = 10$) on the Bray-Curtis dissimilarity of VF profiles. Leiden algorithm was applied with resolution = 1 (see Table S3 for the parameter settings of VERSO employed in the case studies).

Parameter settings of simulations

In order to compare the performance of VERSO step #1 with competing phylogenomic tools (i.e., IQ-TREE¹⁰ and BEAST 2²²), we performed extensive simulations via msprime,⁷⁰ which simulates a backwards-in-time coalescent model.

In particular, we simulated 20 distinct evolutionary processes, with the following parameters: $n = 1,000$ total samples, effective population size $N_e = 0.5$ (i.e., haploid population), mutational rate $M = 2 \times 10^{-6}$ mutations per site per generation, and a genome of length $L = 29,903$ bases. Such parameters were chosen to roughly approximate the mutational rate currently estimated for SARS-CoV-2 (i.e., $M \approx 10^{-3}$ mutations per site per year and $\approx 10^{-3 \frac{\text{generation}}{\text{year}}}$ ¹²²) and to obtain a number of clonal mutations (in the range 15–30) that is comparable with the one observed in the real-world scenarios (see the case studies). As output, msprime returns a phylogenetic tree representing the genealogy between the samples, the genotype of all samples (i.e., the leaves of the tree), and the location of all mutations.

The genotypes of the samples were then inflated with different levels of noise, with false-positive rate α and false-negative rate β (see the parameter settings in Table S1), in order to assess the performance of the methods in conditions of noisy observations and possible sequencing issues. Finally, we subsampled all datasets to obtain two distinct samples sizes (500 and 1,000 samples), in order to test the robustness of methods in conditions of sampling limitations.

The parameters of the phylogenetic methods employed in the comparative assessment are reported in the Supplemental experimental procedures (Table S2).

SUPPLEMENTAL INFORMATION

Supplemental Information can be found online at <https://doi.org/10.1016/j.patter.2021.100212>.

ACKNOWLEDGMENTS

This work was partially supported by the Elixir Italian Chapter and the SysBio-Net project, a Ministero dell'Istruzione, dell'Università e della Ricerca initiative for the Italian Roadmap of European Strategy Forum on Research Infrastructures, and by AIRC-IG grant 22082. Partial support was also provided by the CRUK/AECC/AIRC Accelerator Award #22790, Single-cell Cancer Evolution in the Clinic. We thank Giulio Caravagna, Chiara Damiani, Lucrezia Patrino, and Francesco Craighero for helpful discussions. We also thank David Posada for interesting suggestions on the preliminary version of the manuscript.

AUTHOR CONTRIBUTIONS

D.R., F.A., D.M., A.G., and R.P. designed the approach. D.R., F.A., D.M., and A.G. defined, implemented, and executed the computational methods. D.R., F.A., D.M., and A.G. performed the simulations. D.R., F.A., D.M., C.G.-P., M.A., A.G., and R.P. analyzed the data and interpreted the results. R.P. supervised the experimental data analysis. A.G. and D.R. supervised the computational analysis. A.G. and R.P. drafted the manuscript, which all authors discussed, reviewed, and approved.

DECLARATION OF INTERESTS

The authors declare no competing interests.

Received: June 12, 2020

Revised: November 30, 2020

Accepted: January 22, 2021

Published: January 28, 2021

REFERENCES

- Zhou, P., Yang, X.L., Wang, X.G., Hu, B., Zhang, L., Zhang, W., Si, H.R., Zhu, Y., Li, B., Huang, C.L., et al. (2020). A pneumonia outbreak associated with a new coronavirus of probable bat origin. *Nature* 579, 270–273, <https://doi.org/10.1038/s41586-020-2012-7>.
- Wu, F., Zhao, S., Yu, B., Chen, Y.M., Wang, W., Song, Z.G., Hu, Y., Tao, Z.W., Tian, J.H., Pei, Y.Y., et al. (2020). A new coronavirus associated with human respiratory disease in China. *Nature* 579, 265–269, <https://doi.org/10.1038/s41586-020-2008-3>.
- Andersen, K.G., Rambaut, A., Lipkin, W.I., Holmes, E.C., and Garry, R.F. (2020). The proximal origin of SARS-CoV-2. *Nat. Med.* 26, 450–452, <https://doi.org/10.1038/s41591-020-0820-9>.
- Xiao, K., Zhai, J., Feng, Y., Zhou, N., Zhang, X., Zou, J.J., Li, N., Guo, Y., Li, X., Shen, X., et al. (2020). Isolation of SARS-CoV-2-related coronavirus from Malayan pangolins. *Nature* 583, 286–289, <https://doi.org/10.1038/s41586-020-2313-x>.
- Deng, X., Gu, W., Federman, S., du Plessis, L., Pybus, O.G., Faria, N.R., Wang, C., Yu, G., Bushnell, B., Pan, C.Y., et al. (2020). Genomic surveillance reveals multiple introductions of SARS-CoV-2 into Northern California. *Science* 369, 582–587, <https://doi.org/10.1126/science.abb9263>.
- Grubaugh, N.D., Petrone, M.E., and Holmes, E.C. (2020). We shouldn't worry when a virus mutates during disease outbreaks. *Nat. Microbiol.* 5, 529–530, <https://doi.org/10.1038/s41564-020-0690-4>.
- Grenfell, B.T., Pybus, O.G., Gog, J.R., Wood, J.L., Daly, J.M., Mumford, J.A., and Holmes, E.C. (2004). Unifying the epidemiological and evolutionary dynamics of pathogens. *Science* 303, 327–332, <https://doi.org/10.1126/science.1090727>.
- Poon, L.L., Song, T., Rosenfeld, R., Lin, X., Rogers, M.B., Zhou, B., Sebra, R., Halpin, R.A., Guan, Y., Twaddle, A., et al. (2016). Quantifying influenza virus diversity and transmission in humans. *Nat. Genet.* 48, 195, <https://doi.org/10.1038/ng.3479>.
- Shu, Y., and McCauley, J. (2017). GISAID: global initiative on sharing all influenza data—from vision to reality. *Euro Surveill.* 22, <https://doi.org/10.2807/1560-7917.ES.2017.22.13.30494>.
- Nguyen, L.T., Schmidt, H.A., Von Haeseler, A., and Minh, B.Q. (2015). IQ-TREE: a fast and effective stochastic algorithm for estimating maximum-likelihood phylogenies. *Mol. Biol. Evol.* 32, 268–274, <https://doi.org/10.1093/molbev/msu300>.
- Volz, E.M., Koelle, K., and Bedford, T. (2013). Viral phylodynamics. *PLoS Comput. Biol.* 9, <https://doi.org/10.1371/journal.pcbi.1002947>.
- Faria, N.R., Quick, J., Claro, I., Theze, J., de Jesus, J.G., Giovanetti, M., Kraemer, M.U., Hill, S.C., Black, A., da Costa, A.C., et al. (2017). Establishment and cryptic transmission of Zika virus in Brazil and the Americas. *Nature* 546, 406–410, <https://doi.org/10.1038/nature22401>.
- Worobey, M., Pekar, J., Larsen, B.B., Nelson, M.I., Hill, V., Joy, J.B., Rambaut, A., Suchard, M.A., Wertheim, J.O., and Lemey, P. (2020). The emergence of SARS-CoV-2 in Europe and North America. *Science* 370, 564–570, <https://doi.org/10.1126/science.abc8169>.
- Bandelt, H.J., Forster, P., and Röhl, A. (1999). Median-joining networks for inferring intraspecific phylogenies. *Mol. Biol. Evol.* 16, 37–48, <https://doi.org/10.1093/oxfordjournals.molbev.a026036>.
- O'Meara, B.C. (2012). Evolutionary inferences from phylogenies: a review of methods. *Annu. Rev. Ecol. Evol. Syst.* 43, 267–285, <https://doi.org/10.1146/annurev-ecolsys-110411-160331>.

16. Ronquist, F., Teslenko, M., van der Mark, P., Ayres, D.L., Darling, A., Höhna, S., Larget, B., Liu, L., Suchard, M.A., and Huelsenbeck, J.P. (2012). MrBayes 3.2: efficient Bayesian phylogenetic inference and model choice across a large model space. *Syst. Biol.* *61*, 539–542, <https://doi.org/10.1093/sysbio/sys029>.
17. Stamatakis, A. (2014). RAxML version 8: a tool for phylogenetic analysis and post-analysis of large phylogenies. *Bioinformatics* *30*, 1312–1313, <https://doi.org/10.1093/bioinformatics/btu033>.
18. Didelot, X., Fraser, C., Gardy, J., and Colijn, C. (2017). Genomic infectious disease epidemiology in partially sampled and ongoing outbreaks. *Mol. Biol. Evol.* *34*, 997–1007, <https://doi.org/10.1093/molbev/msw275>.
19. Skums, P., Zelikovsky, A., Singh, R., Gussler, W., Dimitrova, Z., Knyazev, S., Mandric, I., Ramachandran, S., Campo, D., Jha, D., et al. (2018). Quentin: reconstruction of disease transmissions from viral quasispecies genomic data. *Bioinformatics* *34*, 163–170, <https://doi.org/10.1093/bioinformatics/btx402>.
20. De Maio, N., Worby, C.J., Wilson, D.J., and Stoesser, N. (2018). Bayesian reconstruction of transmission within outbreaks using genomic variants. *PLoS Comput. Biol.* *14*, e1006117, <https://doi.org/10.1093/sysbio/sys029>.
21. Kosakovsky Pond, S.L., Weaver, S., Leigh Brown, A.J., and Wertheim, J.O. (2018). HIV-TRACE (TRANsmiSSion Cluster Engine): a tool for large scale molecular epidemiology of HIV-1 and other rapidly evolving pathogens. *Mol. Biol. Evol.* *35*, 1812–1819, <https://doi.org/10.1093/molbev/msy016>.
22. Bouckaert, R., Vaughan, T.G., Barido-Sottani, J., Duchêne, S., Fourment, M., Gavryushkina, A., Heled, J., Jones, G., Kühnert, D., De Maio, N., et al. (2019). Beast 2.5: an advanced software platform for Bayesian evolutionary analysis. *PLoS Comput. Biol.* *15*, e1006650, <https://doi.org/10.1371/journal.pcbi.1006650>.
23. Lai, A., Bergna, A., Acciarri, C., Galli, M., and Zehender, G. (2020). Early phylogenetic estimate of the effective reproduction number of SARS-CoV-2. *J. Med. Virol.* *92*, 675–679, <https://doi.org/10.1002/jmv.25723>.
24. Forster, P., Forster, L., Renfrew, C., and Forster, M. (2020). Phylogenetic network analysis of SARS-CoV-2 genomes. *Proc. Natl. Acad. Sci. U S A* *117*, 9241–9243, <https://doi.org/10.1073/pnas.2004999117>.
25. Dong, R., Pei, S., Yin, C., He, R.L., and Yau, S.S.T. (2020). Analysis of the hosts and transmission paths of SARS-CoV-2 in the COVID-19 outbreak. *Genes* *11*, 637, <https://doi.org/10.3390/genes11060637>.
26. Nakhleh, L. (2009). A metric on the space of reduced phylogenetic networks. *IEEE/ACM Trans. Comput. Biol. Bioinform.* *7*, 218–222, <https://doi.org/10.1126/10.1109/TCBB.2009.2>.
27. Yuan, K., Sakoparnig, T., Markowitz, F., and Beerwinkler, N. (2015). BitPhylogeny: a probabilistic framework for reconstructing intra-tumor phylogenies. *Genome Biol.* *16*, 36, <https://doi.org/10.1186/s13059-015-0592-6>.
28. Villabona-Arenas, C.J., Hanage, W.P., and Tully, D.C. (2020). Phylogenetic interpretation during outbreaks requires caution. *Nat. Microbiol.* *5*, 876–877, <https://doi.org/10.1038/s41564-020-0738-5>.
29. Mavian, C., Marini, S., Manes, C., Capua, I., Prosperi, M., and Salemi, M. (2020). Regaining perspective on SARS-CoV-2 molecular tracing and its implications. *medRxiv*. <https://doi.org/10.1101/2020.03.16.20034470>.
30. Domingo, E., Martínez-Salas, E., Sobrino, F., de la Torre, J.C., Portela, A., Ortín, J., López-Galíndez, C., Pérez-Breña, P., Villanueva, N., Nájera, R., et al. (1985). The quasispecies (extremely heterogeneous) nature of viral RNA genome populations: biological relevance—a review. *Gene* *40*, 1–8, [https://doi.org/10.1016/0378-1119\(85\)90017-4](https://doi.org/10.1016/0378-1119(85)90017-4).
31. Domingo, E., Sheldon, J., and Perales, C. (2012). Viral quasispecies evolution. *Microbiol. Mol. Biol. Rev.* *76*, 159–216, <https://doi.org/10.1128/MMBR.05023-11>.
32. Acevedo, A., Brodsky, L., and Andino, R. (2014). Mutational and fitness landscapes of an RNA virus revealed through population sequencing. *Nature* *505*, 686–690, <https://doi.org/10.1038/nature12861>.
33. Novella, I.S., Domingo, E., and Holland, J.J. (1995). Rapid viral quasispecies evolution: implications for vaccine and drug strategies. *Mol. Med. Today* *1*, 248–253, [https://doi.org/10.1016/s1357-4310\(95\)91551-6](https://doi.org/10.1016/s1357-4310(95)91551-6).
34. Simon-Loriere, E., and Holmes, E.C. (2011). Why do RNA viruses recombine? *Nat. Rev. Microbiol.* *9*, 617–626, <https://doi.org/10.1038/nrmicro2614>.
35. Graudenzi, A., Maspero, D., Angaroni, F., Piazza, R., and Ramazzotti, D. (2020). Mutational signatures and heterogeneous host response revealed via large-scale characterization of SARS-CoV-2 genomic diversity. *iScience* *24*, 102116.
36. Shen, Z., Xiao, Y., Kang, L., Ma, W., Shi, L., Zhang, L., Zhou, Z., Yang, J., Zhong, J., Yang, D., et al. (2020). Genomic diversity of SARS-CoV-2 in coronavirus disease 2019 patients. *Clin. Infect. Dis.* *71*, <https://doi.org/10.1093/cid/ciaa203>.
37. Wölfel, R., Corman, V.M., Guggemos, W., Seilmaier, M., Zange, S., Müller, M.A., Niemeyer, D., Jones, T.C., Vollmar, P., Rothe, C., et al. (2020). Virological assessment of hospitalized patients with COVID-2019. *Nature* *581*, 465–469, <https://doi.org/10.1038/s41586-020-2196-x>.
38. Capobianchi, M.R., Rueca, M., Messina, F., Giombini, E., Carletti, F., Colavita, F., Castilletti, C., Lalle, E., Bordini, L., Vairo, F., et al. (2020). Molecular characterization of SARS-CoV-2 from the first case of COVID-19 in Italy. *Clin. Microbiol. Infect.* *26*, 954–956, <https://doi.org/10.1016/j.cmi.2020.03.025>.
39. Rose, R., Nolan, D.J., Moot, S., Feehan, A., Cross, S., Garcia-Diaz, J., and Lamers, S.L. (2020). Intra-host site-specific polymorphisms of SARS-CoV-2 is consistent across multiple samples and methodologies. *medRxiv*. <https://doi.org/10.1101/2020.04.24.20078691>.
40. Lu, J., du Plessis, L., Liu, Z., Hill, V., Kang, M., Lin, H., Sun, J., François, S., Kraemer, M.U., Faria, N.R., et al. (2020). Genomic epidemiology of SARS-CoV-2 in Guangdong province, China. *Cell* *181*, 997–1003.e9, <https://doi.org/10.1016/j.cell.2020.04.023>.
41. Lythgoe, K.A., Hall, M.D., Ferretti, L., de Cesare, M., MacIntyre-Cockett, G., Trebes, A., Andersson, M., Otecko, N., Wise, E.L., Moore, N., et al. (2020). Shared SARS-CoV-2 diversity suggests localised transmission of minority variants. *bioRxiv*. <https://doi.org/10.1101/2020.05.28.118992>.
42. Seemann, T., Lane, C.R., Sherry, N.L., Duchene, S., Gonçalves da Silva, A., Cally, L., Sait, M., Ballard, S.A., Horan, K., Schultz, M.B., et al. (2020). Tracking the COVID-19 pandemic in Australia using genomics. *Nat. Commun.* *11*, 4376, <https://doi.org/10.1038/s41467-020-18314-x>.
43. Popa, A., Genger, J.W., Nicholson, M.D., Penz, T., Schmid, D., Aberle, S.W., Agerer, B., Lercher, A., Endler, L., Colaço, H., et al. (2020). Genomic epidemiology of superspreading events in Austria reveals mutational dynamics and transmission properties of SARS-CoV-2. *Sci. Transl. Med.* *12*, <https://doi.org/10.1126/scitranslmed.abe2555>.
44. Miralles, R., Gerrish, P.J., Moya, A., and Elena, S.F. (1999). Clonal interference and the evolution of RNA viruses. *Science* *285*, 1745–1747, <https://doi.org/10.1126/science.285.5434.1745>.
45. Xu, D., Zhang, Z., and Wang, F.S. (2004). SARS-associated coronavirus quasispecies in individual patients. *N. Engl. J. Med.* *350*, 1366–1367, <https://doi.org/10.1056/NEJMc032421>.
46. Wright, C.F., Morelli, M.J., Thébaud, G., Knowles, N.J., Herzyk, P., Paton, D.J., Haydon, D.T., and King, D.P. (2011). Beyond the consensus: dissecting within-host viral population diversity of foot-and-mouth disease virus by using next-generation genome sequencing. *J. Virol.* *85*, 2266–2275, <https://doi.org/10.1128/JVI.01396-10>.
47. Park, D., Huh, H.J., Kim, Y.J., Son, D.S., Jeon, H.J., Im, E.H., Kim, J.W., Lee, N.Y., Kang, E.S., Kang, C.I., et al. (2016). Analysis of inpatient heterogeneity uncovers the microevolution of middle east respiratory syndrome coronavirus. *Mol. Case Stud.* *2*, a001214, <https://doi.org/10.1101/mcs.a001214>.
48. Ni, M., Chen, C., Qian, J., Xiao, H.X., Shi, W.F., Luo, Y., Wang, H.Y., Li, Z., Wu, J., Xu, P.S., et al. (2016). Intra-host dynamics of Ebola virus during

2014. *Nat. Microbiol.* 1, 16151, <https://doi.org/10.1038/nmicrobiol.2016.151>.
49. Ramazzotti, D., Caravagna, G., Olde Loohuis, L., Graudenzi, A., Korsunsky, I., Mauri, G., Antoniotto, M., and Mishra, B. (2015). CAPRI: efficient inference of cancer progression models from cross-sectional data. *Bioinformatics* 31, 3016–3026, <https://doi.org/10.1093/bioinformatics/btv296>.
50. Beerenwinkel, N., Schwarz, R.F., Gerstung, M., and Markowetz, F. (2014). Cancer evolution: mathematical models and computational inference. *Syst. Biol.* 64, e1–e25, <https://doi.org/10.1093/sysbio/syu081>.
51. Caravagna, G., Graudenzi, A., Ramazzotti, D., Sanz-Pamplona, R., De Sano, L., Mauri, G., Moreno, V., Antoniotto, M., and Mishra, B. (2016). Algorithmic methods to infer the evolutionary trajectories in cancer progression. *Proc. Natl. Acad. Sci. U S A* 113, E4025–E4034, <https://doi.org/10.1073/pnas.1520213113>.
52. Schwartz, R., and Schäffer, A.A. (2017). The evolution of tumour phylogenetics: principles and practice. *Nat. Rev. Genet.* 18, 213–229, <https://doi.org/10.1038/nrg.2016.170>.
53. Posada, D., and Crandall, K.A. (2001). Intraspecific gene genealogies: trees grafting into networks. *Trends Ecol. Evol.* 16, 37–45, [https://doi.org/10.1016/S0169-5347\(00\)02026-7](https://doi.org/10.1016/S0169-5347(00)02026-7).
54. Boc, A., Diallo, A.B., and Makarenkov, V. (2012). T-rex: a web server for inferring, validating and visualizing phylogenetic trees and networks. *Nucleic Acids Res.* 40 (W1), W573–W579, <https://doi.org/10.1093/nar/cks485>.
55. Bull, J., Badgett, M., Wichman, H.A., Huelsenbeck, J.P., Hillis, D.M., Gulati, A., Ho, C., and Molineux, I. (1997). Exceptional convergent evolution in a virus. *Genetics* 147, 1497–1507. <https://www.genetics.org/content/147/4/1497>.
56. Bandelt, H.J., Quintana-Murci, L., Salas, A., and Macaulay, V. (2002). The fingerprint of phantom mutations in mitochondrial DNA data. *Am. J. Hum. Genet.* 71, 1150–1160, <https://doi.org/10.1086/344397>.
57. Gutierrez, S., Yvon, M., Piroilles, E., Garzo, E., Fereres, A., Michalakos, Y., and Blanc, S. (2012). Circulating virus load determines the size of bottlenecks in viral populations progressing within a host. *PLoS Pathog.* 8, 1–10, <https://doi.org/10.1371/journal.ppat.1003009>.
58. Firestone, S.M., Hayama, Y., Bradhurst, R., Yamamoto, T., Tsutsui, T., and Stevenson, M.A. (2019). Reconstructing foot-and-mouth disease outbreaks: a methods comparison of transmission network models. *Sci. Rep.* 9, 4809, <https://doi.org/10.1038/s41598-019-41103-6>.
59. Wolf, F.A., Angerer, P., and Theis, F.J. (2018). SCANPY: large-scale single-cell gene expression data analysis. *Genome Biol.* 19, 15, <https://doi.org/10.1186/s13059-017-1382-0>.
60. Prosperi, M.C., and Salemi, M. (2012). Qure: software for viral quasispecies reconstruction from next-generation sequencing data. *Bioinformatics* 28, 132–133, <https://doi.org/10.1093/bioinformatics/btr627>.
61. Giallonardo, F.D., Töpfer, A., Rey, M., Prabhakaran, S., Duport, Y., Leemann, C., Schmutz, S., Campbell, N.K., Joos, B., Lecca, M.R., et al. (2014). Full-length haplotype reconstruction to infer the structure of heterogeneous virus populations. *Nucleic Acids Res.* 42, e115, <https://doi.org/10.1093/nar/gku537>.
62. Töpfer, A., Marschall, T., Bull, R.A., Luciani, F., Schönhuth, A., and Beerenwinkel, N. (2014). Viral quasispecies assembly via maximal clique enumeration. *PLoS Comput. Biol.* 10, 1–10, <https://doi.org/10.1371/journal.pcbi.1003515>.
63. Barik, S., Das, S., and Vikalo, H. (2018). QSdpR: viral quasispecies reconstruction via correlation clustering. *Genomics* 110, 375–381, <https://doi.org/10.1016/j.ygeno.2017.12.007>.
64. Knyazev, S., Hughes, L., Skums, P., and Zeilikovsky, A. (2020). Epidemiological data analysis of viral quasispecies in the next-generation sequencing era. *Brief. Bioinformatics*, bbaa101, <https://doi.org/10.1093/bib/bbaa101>.
65. Alizon, S. (2013). Co-infection and super-infection models in evolutionary epidemiology. *Interface Focus* 3, 20130031, <https://doi.org/10.1098/rsfs.2013.0031>.
66. Garcia-Vidal, C., Sanjuan, G., Moreno-García, E., Puerta-Alcalde, P., Garcia-Pouton, N., Chumbita, M., Fernandez-Pittol, M., Pitart, C., Inciarte, A., Bodro, M., et al. (2020). Incidence of co-infections and super-infections in hospitalized patients with COVID-19: a retrospective cohort study. *Clin. Microbiol. Infect.* 27, 83–88, <https://doi.org/10.1016/j.cmi.2020.07.041>.
67. McInnes, L., Healy, J., Saul, N., and Großberger, L. (2018). UMAP: uniform manifold approximation and projection. *J. Open Source Softw.* 3, 861, <https://doi.org/10.21105/joss.00861>.
68. van der Maaten, L., and Hinton, G. (2008). Visualizing data using t-SNE. *J. Machine Learn. Res.* 9, 2579–2605. <http://jmlr.org/papers/v9/vandemaaten08a.html>.
69. Rockett, R.J., Arnott, A., Lam, C., Sadsad, R., Timms, V., Gray, K.A., Eden, J.S., Chang, S., Gall, M., Draper, J., et al. (2020). Revealing COVID-19 transmission in Australia by SARS-CoV-2 genome sequencing and agent-based modeling. *Nat. Med.* 26, 1398–1404, <https://doi.org/10.1038/s41591-020-1000-7>.
70. Kelleher, J., Etheridge, A.M., and McVean, G. (2016). Efficient coalescent simulation and genealogical analysis for large sample sizes. *PLoS Comput. Biol.* 12, 1–22, <https://doi.org/10.1371/journal.pcbi.1004842>.
71. Wakeley, J. (2009). *Coalescent Theory: An Introduction (Roberts and Company Publishers)*.
72. Hadfield, J., Megill, C., Bell, S.M., Huddleston, J., Potter, B., Callender, C., Sagulenko, P., Bedford, T., and Neher, R.A. (2018). Nextstrain: real-time tracking of pathogen evolution. *Bioinformatics* 34, 4121–4123, <https://doi.org/10.1093/bioinformatics/bty407>.
73. Kuhner, M.K., and Felsenstein, J. (1994). A simulation comparison of phylogeny algorithms under equal and unequal evolutionary rates. *Mol. Biol. Evol.* 11, 459–468, <https://doi.org/10.1093/oxfordjournals.molbev.a040126>.
74. Steel, M.A., and Penny, D. (1993). Distributions of tree comparison metrics—some new results. *Syst. Biol.* 42, 126–141, <https://doi.org/10.1093/sysbio/42.2.126>.
75. Bastola, A., Sah, R., Rodriguez-Morales, A.J., Lal, B.K., Jha, R., Ojha, H.C., Shrestha, B., Chu, D.K., Poon, L.L., Costello, A., et al. (2020). The first 2019 novel coronavirus case in Nepal. *Lancet Infect. Dis.* 20, 279–280, [https://doi.org/10.1016/S1473-3099\(20\)30067-0](https://doi.org/10.1016/S1473-3099(20)30067-0).
76. Boni, M.F., Lemey, P., Jiang, X., Lam, T.T.Y., Perry, B.W., Castoe, T.A., Rambaut, A., and Robertson, D.L. (2020). Evolutionary origins of the SARS-CoV-2 sarbecovirus lineage responsible for the covid-19 pandemic. *Nat. Microbiol.* 5, 1408–1417, <https://doi.org/10.1038/s41564-020-0771-4>.
77. Li, X., Giorgi, E.E., Marichanegowda, M.H., Foley, B., Xiao, C., Kong, X.P., Chen, Y., Gnanakaran, S., Korber, B., and Gao, F. (2020). Emergence of SARS-CoV-2 through recombination and strong purifying selection. *Sci. Adv.* 6, eabb9153, <https://doi.org/10.1126/sciadv.abb9153>.
78. Rambaut, A., Holmes, E.C., O’Toole, Á., Hill, V., McCrone, J.T., Ruis, C., du Plessis, L., and Pybus, O.G. (2020). A dynamic nomenclature proposal for SARS-CoV-2 lineages to assist genomic epidemiology. *Nat. Microbiol.* 5, 1403–1407, <https://doi.org/10.1038/s41564-020-0770-5>.
79. O’Toole, A., McCrone, J., and Scher, E. (2020). pangolin 2.0. <https://github.com/cov-lineages/pangolin>.
80. Rambaut, A. (2009). Figtree v1. 3.1. <http://tree.bio.ed.ac.uk/software/figtree/>.
81. Argimón, S., Abudahab, K., Goater, R.J., Fedosejev, A., Bhai, J., Glasner, C., Feil, E.J., Holden, M.T., Yeats, C.A., Grundmann, H., et al. (2016). Microreact: visualizing and sharing data for genomic epidemiology and phylogeography. *Microb. Genomics* 2, e000093, <https://doi.org/10.1099/mgen.0.000093>.

82. Tang, X., Wu, C., Li, X., Song, Y., Yao, X., Wu, X., Duan, Y., Zhang, H., Wang, Y., Qian, Z., et al. (2020). On the origin and continuing evolution of SARS-CoV-2. *Natl. Sci. Rev.* 7, 1012–1023, <https://doi.org/10.1093/nsr/nwaa036>.
83. Zhang, X., Tan, Y., Ling, Y., Lu, G., Liu, F., Yi, Z., Jia, X., Wu, M., Shi, B., Xu, S., et al. (2020). Viral and host factors related to the clinical outcome of covid-19. *Nature* 583, 437–440, <https://doi.org/10.1038/s41586-020-2355-0>.
84. Volz, E.M., Hill, V., McCrone, J.T., Price, A., Jorgensen, D., O’Toole, A., Southgate, J.A., Johnson, R., Jackson, B., Nascimento, F.F., et al. (2020). Evaluating the effects of SARS-CoV-2 spike mutation D614G on transmissibility and pathogenicity. *Cell*. <https://doi.org/10.1016/j.cell.2020.11.020>.
85. Korber, B., Fischer, W.M., Gnanakaran, S., Yoon, H., Theiler, J., Abfalterer, W., Hengartner, N., Giorgi, E.E., Bhattacharya, T., Foley, B., et al. (2020). Tracking changes in SARS-CoV-2 spike: evidence that d614g increases infectivity of the covid-19 virus. *Cell* 182, 812–827, <https://doi.org/10.1016/j.cell.2020.06.043>.
86. Yurkovetskiy, L., Wang, X., Pascal, K.E., Tomkins-Tinch, C., Nyalile, T., Wang, Y., Baum, A., Diehl, W.E., Dauphin, A., Carbone, C., et al. (2020). Structural and functional analysis of the D614G SARS-CoV-2 spike protein variant. *Cell* 183, 739–751, <https://doi.org/10.1016/j.cell.2020.09.032>.
87. Grubaugh, N.D., Hanage, W.P., and Rasmussen, A.L. (2020). Making sense of mutation: what D614G means for the COVID-19 pandemic remains unclear. *Cell* 182, 794–795, <https://doi.org/10.1016/j.cell.2020.06.040>.
88. van Dorp, L., Richard, D., Tan, C.C., Shaw, L.P., Acman, M., and Balloux, F. (2020). No evidence for increased transmissibility from recurrent mutations in SARS-CoV-2. *bioRxiv*. <https://doi.org/10.1101/2020.05.21.108506>.
89. van Dorp, L., Acman, M., Richard, D., Shaw, L.P., Ford, C.E., Ormond, L., Owen, C.J., Pang, J., Tan, C.C., Boshier, F.A., et al. (2020). Emergence of genomic diversity and recurrent mutations in SARS-CoV-2. *Infect. Genet. Evol.* 83, 104351, <https://doi.org/10.1016/j.meegid.2020.104351>.
90. Soares, P., Ermini, L., Thomson, N., Mormina, M., Rito, T., Röhl, A., Salas, A., Oppenheimer, S., Macaulay, V., and Richards, M.B. (2009). Correcting for purifying selection: an improved human mitochondrial molecular clock. *Am. J. Hum. Genet.* 84, 740–759, <https://doi.org/10.1016/j.ajhg.2009.05.001>.
91. Letko, M., Marzi, A., and Munster, V. (2020). Functional assessment of cell entry and receptor usage for RNA and other lineage b betacoronaviruses. *Nat. Microbiol.* 5, 562–569, <https://doi.org/10.1038/s41564-020-0688-y>.
92. Xia, S., Xu, W., Wang, Q., Wang, C., Hua, C., Li, W., Lu, L., and Jiang, S. (2018). Peptide-based membrane fusion inhibitors targeting hcov-229e spike protein HR1 and HR2 domains. *Int. J. Mol. Sci.* 19, 487, <https://doi.org/10.3390/ijms19020487>.
93. Xia, S., Yan, L., Xu, W., Agrawal, A.S., Algaissi, A., Tseng, C.T.K., Wang, Q., Du, L., Tan, W., Wilson, I.A., et al. (2019). A pan-coronavirus fusion inhibitor targeting the hr1 domain of human coronavirus spike. *Sci. Adv.* 5, <https://doi.org/10.1126/sciadv.aav4580>.
94. Li, X., Wang, W., Zhao, X., Zai, J., Zhao, Q., Li, Y., and Chaillon, A. (2020). Transmission dynamics and evolutionary history of 2019-nCoV. *J. Med. Virol.* 92, 501–511, <https://doi.org/10.1002/jmv.25701>.
95. Fay, J.C., Wyckoff, G.J., and Wu, C.I. (2001). Positive and negative selection on the human genome. *Genetics* 158, 1227–1234. <https://www.genetics.org/content/158/3/1227>.
96. Sobel Leonard, A., Weissman, D.B., Greenbaum, B., Ghedin, E., and Koelle, K. (2017). Transmission bottleneck size estimation from pathogen deep-sequencing data, with an application to human influenza a virus. *J. Virol.* 91, <https://doi.org/10.1128/JVI.00171-17>.
97. Ghafari, M., Lumby, C.K., Weissman, D.B., and Illingworth, C.J. (2020). Inferring transmission bottleneck size from viral sequence data using a novel haplotype reconstruction method. *J. Virol.* <https://doi.org/10.1128/JVI.00014-20>.
98. Monsion, B., Froissart, R., Michalakakis, Y., and Blanc, S. (2008). Large bottleneck size in cauliflower mosaic virus populations during host plant colonization. *PLoS Pathog.* 4, 1–7, <https://doi.org/10.1371/journal.ppat.1000174>.
99. Lequime, S., Fontaine, A., Ar Gouilh, M., Moltini-Conclois, I., and Lambrechts, L. (2016). Genetic drift, purifying selection and vector genotype shape dengue virus intra-host genetic diversity in mosquitoes. *PLoS Genet.* 12, 1–24, <https://doi.org/10.1371/journal.pgen.1006111>.
100. Gardy, J.L., and Loman, N.J. (2018). Towards a genomics-informed, real-time, global pathogen surveillance system. *Nat. Rev. Genet.* 19, 9, <https://doi.org/10.1038/nrg.2017.88>.
101. Li, R., Pei, S., Chen, B., Song, Y., Zhang, T., Yang, W., and Shaman, J. (2020). Substantial undocumented infection facilitates the rapid dissemination of novel coronavirus (SARS-CoV-2). *Science* 368, 489–493, <https://doi.org/10.1126/science.abb3221>.
102. Volz, E.M., Pong, S.L.K., Ward, M.J., Brown, A.J.L., and Frost, S.D. (2009). Phylodynamics of infectious disease epidemics. *Genetics* 183, 1421–1430, <https://doi.org/10.1534/genetics.109.106021>.
103. Bedford, T., Greninger, A.L., Roychoudhury, P., Starita, L.M., Famulare, M., Huang, M.L., Nalla, A., Pepper, G., Reinhardt, A., Xie, H., et al. (2020). Cryptic transmission of SARS-CoV-2 in Washington state. *Science* 370, 571–575, <https://doi.org/10.1126/science.abc0523>.
104. Gomez-Carballa, A., Bello, X., Pardo-Seco, J., Martinon-Torres, F., and Salas, A. (2020). Mapping genome variation of SARS-CoV-2 worldwide highlights the impact of COVID-19 super-spreaders. *Genome Res.* <http://www.genome.org/cgi/doi/10.1101/gr.266221.120>.
105. Ramazzotti, D., Angaroni, F., Maspero, D., Ascolani, G., Castiglioni, I., Piazza, R., Antoniotti, M., and Graudenzi, A. (2020). Longitudinal cancer evolution from single cells. *bioRxiv*. <https://doi.org/10.1101/2020.01.14.906453>.
106. Kimura, M. (1969). The number of heterozygous nucleotide sites maintained in a finite population due to steady flux of mutations. *Genetics* 61, 893–903. <https://www.genetics.org/content/61/4/893>.
107. Gusfield, D. (1991). Efficient algorithms for inferring evolutionary trees. *Networks* 21, 19–28, <https://doi.org/10.1002/net.3230210104>.
108. Huson, D.H., and Scornavacca, C. (2012). Dendroscope 3: an interactive tool for rooted phylogenetic trees and networks. *Syst. Biol.* 61, 1061–1067, <https://doi.org/10.1093/sysbio/sys062>.
109. Mindell, D.P. (1993). Merger of taxa and the definition of monophyly (reply to Jan Zrzavý and Zdeněk Skála). *Biosystems* 31, 130–133, [https://doi.org/10.1016/0303-2647\(93\)90041-A](https://doi.org/10.1016/0303-2647(93)90041-A).
110. Zrzavý, J., and Skála, Z. (1993). Holobionts, hybrids, and cladistic classification (reply to David P. Mindell). *Biosystems* 31, 127–130, [https://doi.org/10.1016/0303-2647\(93\)90040-J](https://doi.org/10.1016/0303-2647(93)90040-J).
111. Chan, J.M., Carlsson, G., and Rabadan, R. (2013). Topology of viral evolution. *Proc. Natl. Acad. Sci. U S A* 110, 18566–18571, <https://doi.org/10.1073/pnas.1313480110>.
112. Skála, Z., and Zrzavý, J. (1994). Phylogenetic reticulations and cladistics: discussion of methodological concepts. *Cladistics* 10, 305–313, <https://doi.org/10.1111/j.1096-0031.1994.tb00180.x>.
113. Brandstätter, A., Sängler, T., Lutz-Bonengel, S., Parson, W., Béraud-Colomb, E., Wen, B., Kong, Q.P., Bravi, C.M., and Bandelt, H.J. (2005). Phantom mutation hotspots in human mitochondrial DNA. *Electrophoresis* 26, 3414–3429, <https://doi.org/10.1002/elps.200500307>.
114. Weissensteiner, H., Pacher, D., Kloss-Brandstätter, A., Forer, L., Specht, G., Bandelt, H.J., Kronenberg, F., Salas, A., and Schönherr, S. (2016). HaploGrep 2: mitochondrial haplogroup classification in the era of high-throughput sequencing. *Nucleic Acids Res.* 44 (W1), W58–W63, <https://doi.org/10.1093/nar/gkw233>.
115. Srinivas, G., Möller, S., Wang, J., Künzel, S., Zillikens, D., Baines, J.F., and Ibrahim, S.M. (2013). Genome-wide mapping of gene-microbiota

- interactions in susceptibility to autoimmune skin blistering. *Nat. Commun.* 4, 2462, <https://doi.org/10.1038/ncomms3462>.
116. Pearl, J. (2009). *Causality* (Cambridge University Press).
117. Luecken, M.D., and Theis, F.J. (2019). Current best practices in single-cell RNA-seq analysis: a tutorial. *Mol. Syst. Biol.* 15, <https://doi.org/10.15252/msb.20188746>.
118. Traag, V.A., Waltman, L., and van Eck, N.J. (2019). From Louvain to Leiden: guaranteeing well-connected communities. *Sci. Rep.* 9, 1–12, <https://doi.org/10.1038/s41598-019-41695-z>.
119. Koboldt, D.C., Zhang, Q., Larson, D.E., Shen, D., McLellan, M.D., Lin, L., Miller, C.A., Mardis, E.R., Ding, L., and Wilson, R.K. (2012). VarScan 2: Somatic mutation and copy number alteration discovery in cancer by exome sequencing. *Genome Res.* 22, 568–576, <https://doi.org/10.1101/gr.129684.111>.
120. De Maio, N., Walker, C., Borge, R., Weiguny, L., Slodkowitz, G., and Goldman, N. (2020). Issues with SARS-CoV-2 sequencing data. <https://virological.org/>.
121. Park, D., Tomkins-Tinch, C., Ye, S., Jungreis, I., Shlyakhter, I., Metsky, H., Hanna, Lin, M., Le, V., Lin, A., et al. (2019). broadinstitute/viral-ngs: v1.25.0. <https://doi.org/10.5281/zenodo.3509008>.
122. Bar-On, Y.M., Flamholz, A., Phillips, R., and Milo, R. (2020). Science forum: SARS-CoV-2 (COVID-19) by the numbers. *eLife* 9, e57309, <https://doi.org/10.7554/eLife.57309>.