



HAL
open science

Integrating contact tracing data to enhance outbreak phylogenetic inference: a deep learning approach

Ruopeng Xie, Dillon Adam, Shu Hu, Benjamin Cowling, Olivier Gascuel,
Anna Zhukova, Vijaykrishna Dhanasekaran

► To cite this version:

Ruopeng Xie, Dillon Adam, Shu Hu, Benjamin Cowling, Olivier Gascuel, et al.. Integrating contact tracing data to enhance outbreak phylogenetic inference: a deep learning approach. *Molecular Biology and Evolution*, In press, pp.msae232. 10.1093/molbev/msae232 . pasteur-04768262

HAL Id: pasteur-04768262

<https://pasteur.hal.science/pasteur-04768262v1>

Submitted on 5 Nov 2024

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution - NonCommercial 4.0 International License

Integrating contact tracing data to enhance outbreak phylogenetic inference: a deep learning approach

Ruopeng Xie^{1,2,*}, Dillon C. Adam¹, Shu Hu^{1,2}, Benjamin J. Cowling^{1,3}, Olivier Gascuel⁴, Anna Zhukova^{5,6*},
Vijaykrishna Dhanasekaran^{1,2,*}

Affiliations:

¹School of Public Health, LKS Faculty of Medicine, The University of Hong Kong; Hong Kong S.A.R., China.

²HKU-Pasteur Research Pole, School of Public Health, LKS Faculty of Medicine, The University of Hong Kong; Hong Kong S.A.R., China

³Laboratory of Data Discovery for Health, Hong Kong Science and Technology Park, New Territories, Hong Kong S.A.R., China.

⁴Institut de Systématique, Evolution, Biodiversité (ISYEB, UMR 7205 – CNRS, MNHN, SU, EPHE, UA), Muséum National d'Histoire Naturelle, 45 rue Buffon, 75005 - Paris, France.

⁵Bioinformatics and Biostatistics Hub, Institut Pasteur, Université de Paris, 75015 Paris, France.

⁶G5 Evolutionary Dynamics of Infectious Diseases, Institut Pasteur, Université de Paris, 75015 Paris, France

*Corresponding authors. Email: rxie@connect.hku.hk, anna.zhukova@pasteur.fr, veej@hku.hk

Abstract

Phylogenetics is central to understanding infectious disease dynamics through the integration of genomic and epidemiological data. Despite advancements, including the application of deep learning to overcome computational limitations, significant challenges persist due to data inadequacies and statistical unidentifiability of key parameters. These issues are particularly pronounced in poorly resolved phylogenies, commonly observed in outbreaks such as SARS-CoV-2. In this study, we conducted a thorough evaluation of PhyloDeep, a deep learning inference tool for phylogenetics, assessing its performance on poorly resolved phylogenies. Our findings reveal the limited predictive accuracy of PhyloDeep (and other state-of-the-art approaches) in these scenarios. However, models trained on poorly resolved, realistically simulated trees

© The Author(s) 2024. Published by Oxford University Press on behalf of Society for Molecular Biology and Evolution. This is an Open Access article distributed under the terms of the Creative Commons Attribution-NonCommercial License (<https://creativecommons.org/licenses/by-nc/4.0/>), which permits non-commercial re-use, distribution, and reproduction in any medium, provided the original work is properly cited. For commercial re-use, please contact reprints@oup.com for reprints and translation rights for reprints. All other permissions can be obtained through our RightsLink service via the Permissions link on the article page on our site—for further information please contact journals.permissions@oup.com.

1 demonstrate improved predictive power, despite not being infallible, especially in scenarios with
2 superspreading dynamics, whose parameters are challenging to capture accurately. Notably, we observe
3 markedly improved performance through the integration of minimal contact tracing data, which refines
4 poorly resolved trees. Applying this approach to a sample of SARS-CoV-2 sequences partially matched to
5 contact tracing from Hong Kong yields informative estimates of superspreading potential, extending beyond
6 the scope of contact tracing data alone. Our findings demonstrate the potential for enhancing phylodynamic
7 analysis through complementary data integration, ultimately increasing the precision of epidemiological
8 predictions crucial for public health decision making and outbreak control.

10 **Introduction**

11 Phylogenetic analysis of genomic sequence data offers a powerful toolkit for understanding the emergence,
12 spread, and evolution of infectious diseases. As an interdisciplinary field, phylodynamics aims to integrate
13 genomic and epidemiological data in a unified framework to extract detailed insights into epidemic history
14 (Drummond et al., 2005; Stadler et al., 2013; Volz et al., 2009), population dynamics (Stadler & Bonhoeffer,
15 2013; Volz et al., 2009), and disease emergence (Pekar et al., 2022; Worobey et al., 2014). Its key
16 advantage lies in providing independent information regarding epidemic history, complementing traditional
17 epidemiological surveillance data (Vaughan et al., 2024; Voznica et al., 2022). This makes it invaluable for
18 validating and substantiating findings from epidemiological modelling, particularly in contexts where
19 conventional surveillance data are scarce and genomic sampling is randomized.

20 However, many conventional phylodynamic models based on likelihood approaches (e.g. maximum
21 likelihood estimation and Bayesian approaches) are computationally intensive and can become practically
22 unfeasible as the number of taxa increases (Hohna & Drummond, 2012). Addressing this issue sometimes
23 involves likelihood-free methods such as approximate Bayesian computation (ABC) (Saulnier et al., 2017),
24 which sidestep the need for direct likelihood calculations. More recently, deep learning methods such as
25 PhyloDeep (Voznica et al., 2022) have emerged as another potential solution, enabling rapid estimation of
26 epidemiological parameters from large phylogenetic trees in a matter of seconds. To achieve this, PhyloDeep
27 utilizes deep neural network models trained against phylogenies simulated under well-established birth-
28 death models: the basic birth-death model (BD) (Leventhal et al., 2014; Stadler et al., 2012), the birth-death
29 model with exposed and infectious classes (BDEI) (Kuhnert et al., 2016; Stadler et al., 2013), and the birth-
30 death model with superspreading (BDSS) (Stadler et al., 2013). PhyloDeep has also been validated for
31 diversification analyses (Lambert et al., 2023) and viral phylogeography (Thompson et al., 2024).

32 Despite these methodological advancements, critical challenges remain concerning the adequacy of datasets
33 and the statistical identifiability of the parameters of interest from sequence data. This issue is particularly
34 pronounced for viral sequences arising from epidemics and outbreaks, which frequently yield many identical
35 sequences, resulting in poorly resolved phylogenies with numerous polytomies. Examples include SARS-CoV-

1 2, Mpox (monkeypox) virus (Paredes et al., 2024), and Respiratory syncytial virus (RSV) (Eden et al., 2022).
2 These poorly resolved trees typically do not align with the “idealistic”, well-resolved trees posited by
3 phylodynamic models like birth-death models, where branching events are assumed to correspond to
4 transmission events. Such misalignment could introduce biases, compromising the accuracy and reliability of
5 inference methods and potentially leading to incorrect interpretations of epidemic dynamics and disease
6 transmission.

7 To address these concerns, this study utilizes the PhyloDeep framework to assess the impact of potential
8 biases introduced by poorly resolved phylogenies, using the SARS-CoV-2 as an example of a virus outbreak
9 characterized by the BDSS model, which splits population into normal and superspreaders while tracking
10 superspreading potential (**Fig. 1**). Our analysis reveals that neural network models in PhyloDeep (and other
11 state-of-the-art approaches) struggle to precisely predict epidemiological parameters when applied to
12 poorly resolved phylogenetic trees, but performance does improve when models are trained on poorly
13 resolved, realistically simulated phylogenies rather than on “idealistic” trees from birth-death models, as
14 previously done in PhyloDeep. However, capturing superspreading dynamics remains a challenge. Notably,
15 integrating contact tracing data substantially enhances predictive accuracy by constraining tree space and
16 aligning them more closely with “idealistic” trees. Additionally, this integration also proves beneficial in the
17 Bayesian inference framework implemented in BEAST2 (Bouckaert et al., 2014). We illustrate these findings
18 using real SARS-CoV-2 data collected during the third and fourth waves of the epidemic in Hong Kong.

20 **Results**

21 Building on the PhyloDeep approach, we simulated phylogenetic trees (idealistic) using the BDSS model,
22 covering a broad range of epidemiological parameter values associated with the SARS-CoV-2 virus (**Fig.**
23 **1**). These simulated trees were transformed into six additional forms, ranging from idealized simulations to
24 those reflecting the complexities of real-world sequences and trees (**Figs. 1 and 2**). Neural networks trained
25 with summary statistics (SSs) were applied to each tree type to perform regression tasks, estimating
26 epidemiological parameters and evaluating the performance of these models comprehensively.

28 **Simulations of phylogenetic trees**

29 Initially, we simulated 200,000 time-scaled trees using the BDSS model (**Fig. 2, baseline tree**). These trees
30 serve as our reference “idealistic” trees and capture transmission events at internal nodes consistent with the
31 PhyloDeep framework. To emulate SARS-CoV-2 phylogenetic trees, all baseline trees were transformed into
32 genetic distance trees (**Fig. 2, genetic baseline tree**). This transformation relied on a binomial distribution of
33 mutation counts given a mean substitution rate of 8×10^{-4} per site per year, resulting in approximately 24
34 mutations observed annually for a sequence length of 29903 bases (see methods for details). Branches with

1 lengths representing zero mutation were collapsed, resulting in trees with polytomies (**Fig. 2, genetic**
 2 **polytomous tree**), which were then randomly resolved using a coalescent approach, yielding binary trees
 3 (**Fig. 2, genetic resolved tree**). The number and size of polytomies in our simulated trees varied from 1 to
 4 170 and 3 to 934, respectively, with a total tip range of 200 to 1000, encompassing those observed in
 5 SARS-CoV-2 trees in Hong Kong (Supplementary Figure S1). Lastly, each of the three transformed genetic
 6 distance trees were dated using LSD2 (To et al., 2016) (**Fig. 2, dated baseline tree, dated polytomous tree,**
 7 **dated resolved tree**). Genetic Polytomous Trees, Genetic Resolved Trees, Dated Polytomous Trees, and
 8 Dated Resolved Trees, represent entirely altered topologies and are deemed poorly resolved, realistic trees.
 9 They are analogous to trees inferred from sequencing data using established software such as RAxML-NG
 10 (Kozlov et al., 2019), IQ-TREE (Nguyen et al., 2015), PhyML (Guindon et al., 2010), FastTree (Price et al.,
 11 2010) or TreeTime (Sagulenko et al., 2018). In contrast, the remaining three types, namely Baseline Trees,
 12 Genetic Baseline Trees, and Dated Baseline Trees, retain a known correct topology that cannot be derived
 13 from sequence data alone (**Fig. 2**).

14

15 **Performance comparison of neural network models for each type of phylogenetic tree**

16 We utilized a dataset totalling 199,000 trees to train the neural network models, reserving 1,000 trees for
 17 validation purposes. Ensuring consistency across the models, we utilized the same 99 summary statistics (SSs)
 18 representation and feed-forward neural networks (FFNNs) architectures for each tree type, as used in
 19 PhyloDeep (**Fig. 1**). Specifically, for the three types of genetic distance trees, namely Genetic Baseline Trees,
 20 Genetic Polytomous Trees and Genetic Resolved Trees, we adapted the 99 SSs designed for time-scaled
 21 trees to 90 SSs for genetic distance trees (refer to the Methods section). Consequently, we trained seven
 22 neural network models: Baseline-Model, Dated Baseline-Model, Dated Resolved-Model, Dated Polytomous-
 23 Model, Genetic Baseline-Model, Genetic Resolved-Model, and Genetic Polytomous-Model.

24 Our results show that models trained and tested on trees with unchanged topologies (i.e. Baseline-Model,
 25 Dated Baseline-Model, and Genetic Baseline-Model) did well in predicting all parameters. Estimates for R_0
 26 and $1/\gamma$ tended to exhibit greater accuracy compared to superspreading parameters (X_{ss} and f_{ss}) (**Fig. 3A**
 27 and Supplementary Table S2), which is consistent with the findings from PhyloDeep (Voznica et al., 2022).
 28 As expected, the Baseline-Model exhibited the best performance, achieving mean relative errors of 0.095
 29 for R_0 , 0.092 in $1/\gamma$, 0.215 for X_{ss} and 0.167 for f_{ss} . Conversely, models trained and tested on trees with
 30 altered topologies (Dated Resolved-Model, Dated Polytomous-Model, Genetic Polytomous-Model and
 31 Genetic Resolved-Model) encountered challenges in accurately predicting superspreading parameters. This
 32 suggests that phylogenetic trees with polytomies lack sufficient phylogenetic resolution to accurately recover
 33 parameters related to superspreading. Models trained and tested on dated trees generally outperformed
 34 those trained and tested on the equivalent genetic distance trees in most scenarios, demonstrating the value
 35 of tip dates for informing model learning and estimating parameters.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34

Impact of poorly resolved phylogenetic trees on models trained with “idealistic” trees

To evaluate the impact of using poorly resolved realistic phylogenetic trees as input on neural network models trained with “idealistic” trees, we tested the Baseline-Model and Dated Baseline-Model with 1,000 Dated Resolved Trees and the Genetic Baseline-Model with 1,000 Genetic Resolved Trees (Fig. 3 and Supplementary Table S2). The results revealed that the relative error for each parameter was approximately twice as high or more compared to when using “idealistic” test trees (Fig. 3B). Notably, the relative errors for the superspreading parameters (X_{ss} and f_{ss}) were around or exceeded 0.5 (50%). This demonstrates that models trained on “idealistic” trees struggled to predict accurately epidemiological parameters from poorly resolved, realistic phylogenetic trees. Conversely, models trained on poorly resolved trees (such as Genetic Polytomous, Genetic Resolved, Dated Polytomous, and Dated Resolved) performed better, underscoring the importance of training on data that mirrors real-world complexity (Fig. 3A). However, despite improvements, the higher predictive errors specific to superspreading parameters relative to other epidemiological parameters seemed to persist (Fig. 3), highlighting the inherent challenge in estimating superspreading potential from such poorly resolved trees. Additionally, despite repeatedly generating different Genetic Resolved and Dated Resolved trees from the polytomous trees as input, the predicted parameters tended to converge towards similar estimates, which differed substantially from the actual parameters originally input, thus indicating a form of bias in the estimations.

Improving predictions by integrating contact tracing data

To improve model accuracy, a reasonable approach involves correcting the observed topology of input trees so that they closely resemble the equivalent “idealistic” trees. In this context, we investigated the potential of leveraging contact tracing data (including cluster information and infection times) to aid in refining the topology of Genetic Polytomous trees, for example, to match Baseline or Dated Baseline trees to varying extents (Supplementary Figure S2). We derived contact tracing information from the simulated Baseline trees, treating all descendants of each internal node as a cluster, with the dates of internal nodes considered as infection times of each cluster’s index case (Supplementary Figure S3). With this addition of cluster information and assuming perfect observation, the topology of Genetic Polytomous trees can be fully corrected (matching the genetic baseline trees), with external nodes subsequently dated to produce Dated Baseline trees (Supplementary Figure S2). Furthermore, if the infection times of clusters are known, time constraints can also be applied to internal nodes, effectively recovering equivalent Baseline trees from the Genetic Polytomous trees. In real-world scenarios, however, the extent of case observation is often limited and imperfect, and the accuracy of any available contact tracing data is uncertain and subject to additional biases.

1 Therefore, to assess how the quantity of contact tracing data influences our predictions within the context of
2 phylogenetic trees, we simulated scenarios where 0%, 25%, 50%, 75%, and 100% of internal nodes were
3 randomly selected to provide cluster information and infection times. We then evaluated the performance
4 of the Baseline-Model and Dated Baseline-Model (**Fig. 4A** and Supplementary Table S3). The former
5 requires cluster information to resolve polytomies and infection times, with a time constraint margin of 1 day,
6 to estimate the lengths of newly created internal branches from Genetic Polytomous Trees (Supplementary
7 Figure S2G), while the latter relies solely on cluster information (Supplementary Figure S2F). For any
8 remaining nodes lacking contact tracing data, we resolved them randomly as before. Our results indicated
9 that even with just 25% of contact tracing data incorporated, the mean relative errors for R_0 and $1/\gamma$ could
10 be reduced to below 0.2, representing an improvement of 48% to 66% (Supplementary Table S3). As the
11 availability of contact tracing data increased, model performance consistently improved, particularly in
12 predicting superspreading parameters as could be expected. Incorporating 50% or more of contact tracing
13 data yielded estimates of superspreading parameters, with mean relative errors around or below 30%,
14 achieving an improvement of at least 22% (Supplementary Table S3). Notably, the Dated Baseline-Model
15 generally outperformed the Baseline-Model except when contact tracing was 100% complete and a harsh
16 time constraint margin of 0.1 day (Supplementary Table S3). Furthermore, the Dated Baseline-Model only
17 required cluster information to refine the input trees, suggesting its greater relevance to real-world scenarios.

18 We compared the Dated Baseline-Model and Baseline-Model to the gold standard likelihood-based
19 Bayesian tool BEAST2 (Bouckaert et al., 2014), across varying levels of contact tracing data. BEAST2's
20 performance improved with increased proportions of contact tracing data, which includes cluster information
21 and infection times (**Fig. 4B** and Supplementary Table S4). However, BEAST2 consistently underperformed
22 compared to our Baseline-Model, except when no contact tracing data was incorporated. Even in this
23 scenario, it still performed worse than models trained on poorly resolved phylogenies (**Fig. 3A** and
24 Supplementary Tables S2 and S4). Additionally, BEAST2 struggled to accurately infer superspreading
25 parameters, even with 100% contact tracing data, which aligns with the findings of PhyloDeep (Voznica et
26 al., 2022). Further, providing only cluster information, which modifies the tree topology without correcting
27 the time of internal nodes, did not substantially enhance BEAST2's performance (**Fig. 4C** and Supplementary
28 Table S4), likely due to the loss of this crucial temporal information.

29

30 **Case study of SARS-CoV-2 waves in Hong Kong**

31 By 2022, Hong Kong had effectively controlled the local spread of SARS-CoV-2, experiencing four
32 significant waves during which extensive sequence sampling and epidemiological surveillance were
33 conducted, as detailed in our previous study (Gu et al., 2022). To demonstrate our method of integrating
34 contact tracing data to improve model prediction, we used real-world SARS-CoV-2 data from the third and
35 fourth waves in Hong Kong, analysed from May 13 to August 1, 2020 (460 sequences and 1,930 local

1 cases), and from September 30 to December 8, 2020 (243 sequences and 1,577 local cases). Utilizing all
2 available SARS-CoV-2 sequences from these periods along with partial contact tracing data (only cluster
3 information available), covering 16.56 % for the third wave and 9.50% for the fourth wave (see Methods),
4 we evaluate the differences in prediction outcomes when using the Dated Baseline-Model, with input trees
5 refined by contact tracing data (Dated Resolved-Cluster) and without it (Dated Resolved, random resolution
6 of polytomies).

7 Initially, we verified the suitability of the input trees generated by RAxML-NG (Kozlov et al., 2019) using
8 the GTR+G4+FO substitution model with random resolution of polytomies, through principal component
9 analysis (PCA) and by comparing the range of each simulated SS to ensure the models and scenarios were
10 predictive. All trees from Hong Kong passed this PCA check, but seven SSs related to transmission chain
11 features for the Dated Resolved tree of wave 4 were outside the [min, max] range of the simulated values
12 (Supplementary Figure S4 and Table S5). After integrating the available contact tracing data (9.50%, as
13 detailed in the Methods), only one SS remained outside the simulated range, albeit very close to the lower
14 boundary (Supplementary Table S5).

15 The prediction results indicated a notable change when contact tracing data was used to refine tree
16 topology, especially for wave 4 (**Table 1**). With the Dated Resolved-Cluster tree, we estimated an R_0 of
17 1.59 and 1.52, infection-to-sampling periods (infectious periods, $1/\gamma$) of 4.6 and 8.6 days, X_{ss} of 8.1 and
18 16.4, f_{ss} of 0.091 and 0.078 for waves 3 and 4, respectively. Given X_{ss} and f_{ss} , we can calculate the
19 dispersion value k (see Methods), which is commonly used as a measure of superspreading potential. For
20 waves 3 and 4 we calculated $k = 0.47$ and 0.25 respectively, where lower values of k represent increasing
21 superspreading potential. Conversely, using the Dated Resolved tree, we estimated an R_0 of 1.70 and 2.06,
22 infection-to-sampling periods of 5.7 and 20.1 days, X_{ss} of 7.6 and 7.2, f_{ss} of 0.090 and 0.076, and k of
23 0.49 and 0.66 for waves 3 and 4, respectively. The unusually long infection-to-sampling periods of 20.1
24 days observed in wave 4 may be attributed to the seven SSs that exceeded the expected range, which
25 likely influenced these skewed predictions (Supplementary Figure S4 and Table S5). Further, based solely
26 on epidemiological records, we estimated an R_0 of 1.69 and 1.93, and k of 0.45 and 0.26 for waves 3 and
27 4, separately (**Table 1**). The observed discrepancies highlight the critical need for integrating diverse data
28 sources and analytical methods in estimating epidemiological parameters, thereby enabling a more
29 comprehensive and systematic understanding of epidemic dynamics.

30 Additionally, we conducted 200 random resolution of polytomies for these SARS-CoV-2 trees to measure
31 the robustness of the predictions. The resulting standard deviation were notably small (**Table 1**), indicating
32 that the predictions were not significantly affected by the random resolution of polytomies, suggesting our
33 models could efficiently extract essential cluster information and guide robust predictions. The 95%
34 confidence intervals (CIs) were generated by parametric bootstrap as per the methodology of PhyloDeep.

1 The substantial width of CIs for superspreading parameters again highlight the inherent difficulty in
2 predicting these metrics.

3
4 **Table 1.** Comparison of inference of epidemiological parameters based on waves 3 and 4 of SARS-CoV-2
5 in Hong Kong.

| Waves | Input tree | R_0 | Infection-to-sampling period (day) | X_{ss} | f_{ss} | Dispersion k |
|-------|----------------------------|-------------------------------|------------------------------------|---------------------------------|-------------------------------|-------------------------|
| 3 | Dated Resolved | 1.699±0.096 (1.460, 2.172) | 5.720±1.018 (4.427, 10.804) | 7.608±1.496 (4.141, 18.696) | 0.090±0.022 (0.057, 0.163) | 0.488 (0.441, 0.543) |
| | Dated Resolved-Cluster | 1.588±0.077 (1.330, 1.993) | 4.636±0.635 (3.373, 8.238) | 8.078±1.709 (3.911, 17.733) | 0.091±0.021 (0.054, 0.167) | 0.467 (0.418, 0.517) |
| | Epidemiological inference* | 1.693 (1.649, 1.738) | NA | NA | NA | 0.451 (0.421, 0.481) |
| 4 | Dated Resolved | 2.062±0.072 (1.628, 3.220) | 20.071±1.663 (14.235, 32.668) | 7.232±1.423 (2.197, 23.198) | 0.076±0.009 (0.050, 0.154) | 0.658 (0.596, 0.737) |
| | Dated Resolved-Cluster | 1.518±0.091 (1.284, 2.055) | 8.629±0.881 (6.548, 14.929) | 16.388±2.692 (5.895, 33.409) | 0.078±0.007 (0.050, 0.161) | 0.250 (0.227, 0.278) |
| | Epidemiological inference* | 1.933 (1.858, 2.012) | NA | NA | NA | 0.264 (0.248, 0.279) |

6 Note: Values predicted by neural network models are expressed as mean ± standard deviation generated by randomly
7 resolving polytomies $n = 200$ times. Values in parentheses are the 95% CI. In the BDSS model, the term “infectious period”
8 refers to the interval from the time of infection to the sampling date. To prevent confusion in epidemiological contexts, we
9 have opted to use “infection-to-sampling period” in place of “infectious period”. * Epidemiological inference uses a
10 combination of line-listed incidence data to estimate R_0 and contact tracing data to estimate k .

11 12 13 Discussion

14 In this study, we assessed the performance of established neural network models (PhyloDeep) in predicting
15 epidemiological parameters and the applicability of these models to real-world scenarios using SARS-CoV-
16 2 as a case study for both simulation and empirical analyses. Our findings demonstrate the relative
17 performance limitations of utilizing neural network models trained on simulated phylogenetic trees
18 (“idealistic” trees) when predicting parameters from poorly resolved trees (“realistic” trees), and show that
19 models alternatively trained on simulated trees of similar resolution can improve the accuracy of predictions.

1 Beyond upstream improvements to model training, we show that by using contact tracing data to partially
2 resolve the topology and node dates of input trees downstream, additional performance enhancements can
3 be achieved. We apply this approach to SARS-CoV-2 genome sequences from Hong Kong matched to
4 minimal contact tracing data, producing new phylodynamic estimates of both R_0 (basic reproductive number)
5 and k (dispersion measure of superspreading potential).

6 Without the incorporation of contact tracing data, we found that our improved models trained on simulated
7 poorly resolved trees still struggled to accurately estimate parameters related to superspreading, even when
8 attempting to overfit neural network models on smaller subsets of trees (Supplementary Table S6). This issue
9 is particularly pronounced when sequences are nearly identical, like for SARS-CoV-2, which results in
10 potentially biased estimations likely to misinform public health decision makers. Traditional phylodynamic
11 inference methods (e.g. maximum likelihood estimation and Bayesian approaches) with models that assume
12 ideal binary trees and not representing sequence evolution, also struggle in parameter estimation under
13 these conditions (Supplementary Table S4) (Lewis et al., 2005; Morel et al., 2021). Together this emphasizes
14 the importance of incorporating even minimal contact tracing data as we have done in our study, but also
15 utilizing more comprehensive summary statistics focused on clusters and polytomies that can effectively
16 capture the complexity of the underlying transmission dynamic. One previous study (Tran-Kiem & Bedford,
17 2024) has demonstrated a connection between the size distribution of identical sequence clusters and
18 transmission dynamics, however, our attempts to incorporate similar information into our neural network
19 models, trained on genetic distance trees, yielded limited improvements. As an ongoing area of research
20 interest, future studies could evaluate the relative predictive performance of models that expand the
21 potential range of summary statistics related to clusters and polytomies, and experiment with alternate
22 architectures such as Graph Neural Networks (GNN) and Convolutional Neural Networks (CNN)
23 incorporating a more complete representations of trees, such as Compact Bijective Ladderized Vectors
24 (CBLV) (Voznica et al., 2022).

25 Besides superspreading, the incubation period is another significant aspect of pathogen transmission
26 dynamics. For example, estimates of the SARS-CoV-2 incubation period were used to justify the World
27 Health Organization's (WHO) recommendation of a 14-day quarantine period for contacts of infected cases
28 (Wells et al., 2021). In our approach, we utilized a BDSS model, which does not account for the incubation
29 period, but defines the infectious period as the interval from infection time to sampling date otherwise known
30 as the delay interval. Employing the Dated Baseline-Model with the Dated Resolved-Cluster tree, we
31 determined the infectious period/delay interval of waves 3 and 4 to be approximately one week, however
32 the delay for wave 4 was longer than that for wave 3, suggesting case detection speed was somewhat
33 challenged. The longer delay in wave 4 could be explained by the sudden rise in cases associated with the
34 largest single SARS-CoV-2 superspreading event detected in Hong Kong prior to widespread vaccination,
35 which also triggered the start of wave 4 (Adam et al., 2022; Gu et al., 2022).

1 Remarkably, the estimation of R_0 exhibited robust performance across our neural network models, with
2 models trained on dated trees outperforming those based on genetic distance trees. This underscores the
3 value of tip dates for R_0 estimation, particularly as sequence variability decreases. This is in line with recent
4 studies that highlight the increasing importance of sampling dates for phylodynamic inference when sequence
5 variability is low (Featherstone et al., 2023). When poorly resolved trees were used as input, models like
6 the Dated Resolved-Model and Dated Polytomous-Model showed excellent performance, suggesting their
7 potential for effective and accurate R_0 and $1/\gamma$ predictions from sequence data. This offers a promising
8 avenue for tracking epidemic dynamics using sequence data, which, when compared with epidemiological
9 records, can provide deeper insights and mitigate potential sampling biases. Future investigations are
10 needed to ascertain the extent to which sequence data can facilitate robust predictions and to evaluate the
11 effects of progressively incorporating new sequence samples.

12 Our study acknowledges certain limitations. Notably, the BDSS model does not account for the incubation
13 period of the disease, introducing a significant source of uncertainty. The omission of the incubation period
14 from our transmission models necessitates further exploration in future studies to mitigate these uncertainties.
15 For example, an alternative approach could use a Susceptible-Exposed-Infected-Recovered (SEIR) model
16 with a superspreading compartment, grounded in structured coalescent theory (Volz & Siveroni, 2018), which
17 has been used to study superspreading and nonlinear incidence in SARS-CoV-2 studies (Geidelberg et al.,
18 2021; Miller et al., 2020; Moreno et al., 2020; Ragonnet-Cronin et al., 2021). Additionally, real-world
19 contact tracing data may contain inherent biases and inaccuracies. In applying our model to the SARS-CoV-
20 2 dataset from Hong Kong, we presumed the accuracy of the contact tracing data. This assumption allowed
21 us to collapse all associated children (see Methods), including those are not recorded within the cluster,
22 potentially leading to an inaccurate refinement of the tree topology and biased predictions. Our primary
23 epidemiological inference of R_0 assumed a comparable SIR model of transmission and an exponentially
24 distributed generation time like BDSS, though tended to be slightly higher than the mean R_0 estimated from
25 PhyloDeep (Table 1). This method, which links the initial growth rate of an epidemic to R_0 (Wallinga &
26 Lipsitch, 2007) is however known to exhibit a slight upward bias for smaller R_0 values ($R_0 < 2$). (Obadia et
27 al., 2012). Further sensitivity analyses assuming gamma-distributed generation times, unlike BDSS, resulted
28 in even higher values R_0 , partially validating the results from our Dated Baseline-Model with Dated Resolved-
29 Cluster tree (Supplementary Table S8).

30 Importantly, making trees poorly resolved during training hinges on the specific sequence length and
31 evolution rate of SARS-CoV-2, rendering the neural networks trained in this study inapplicable to other
32 viruses. To extend their use to other pathogens, modifications are required to accommodate variations in
33 sequence length and evolution rate, training pathogen-specific neural networks as we show for SARS-CoV-
34 2. This contrasts with PhyloDeep, which was designed for studying a diverse array of pathogens.

1 Correspondingly, the choice of a specific birth-death model emerges as another crucial factor that must be
2 carefully considered.

3 Overall, this study highlights the challenges of relying solely on viral phylogenetic trees generated from
4 sequences for estimating superspreading events. The integration of even minimal contact tracing data can
5 significantly enhance model predictions, emphasizing the importance of such data in surveillance efforts for
6 emerging infectious diseases, particularly when viral sequences lack variability. We hope our comprehensive
7 evaluation will not only enhance deep learning applications but also extend beyond, enriching established
8 methodologies within phylogenetics and phylodynamics.

9

10 **Methods**

11 **Simulations**

12 In this study, SARS-CoV-2 served as the reference pathogen for evaluating the performance of the existing
13 deep learning model PhyloDeep. Given the marked overdispersion in SARS-CoV-2 transmission dynamics,
14 characterized by superspreading (Adam et al., 2020; Du et al., 2022; Guo et al., 2022), we used
15 treesimulator (v0.1.7: (Zhukova & Gascuel, 2024)) to generate time-scaled phylogenetic trees (detailed in
16 Supplementary Table S1). These trees were generated with a BDSS model, distinguishing cases into
17 superspreaders (S) and normal spreaders (N), in addition to the conventional parameterization of the Birth-
18 Death model, i.e. R_0 and $1/\gamma$. Superspreaders constitute a small fraction of the total simulated population
19 (denoted by $f_{SS} = \beta_{SS}/(\beta_{SS} + \beta_{SN})$) but can transmit the virus at rates significantly higher than normal
20 spreaders, where the superspreading transmission ratio is denoted as $X_{SS} = \beta_{SS}/\beta_{NS} = \beta_{SN}/\beta_{NN}$. Upon
21 reviewing the 98 summary statistics (SS) (see details in Feature representation and neural network models
22 section), it was noted that certain metrics associated with branch lengths and superspreading events based
23 on the SARS-CoV-2 dataset from Hong Kong fell outside the [min, max] range of simulated values in
24 PhyloDeep, characterized by a lower median/mean SS and increased variance SS (detailed in
25 Supplementary Table S7). Consequently, to better capture the complexities of SARS-CoV-2 transmission
26 dynamics, we expanded the range of epidemiological parameters for tree simulation in PhyloDeep,
27 summarized in Supplementary Table S1.

28 Simulated time-scaled trees are transformed into Genetic Baseline trees, with branch lengths determined by
29 a binomial process, B (n =sequence length, p =evolutionary rate \times branch length of time-scaled trees). For
30 SARS-CoV-2, the sequence length is 29,903, and the evolutionary rate has a mean of 8×10^{-4} and a standard
31 deviation of 4×10^{-4} substitutions per site per year, with a lognormal distribution (Hadfield et al., 2018; Jolly
32 & Scaria, 2021). In Genetic Baseline trees, branches representing zero mutation are collapsed to form
33 Genetic Polytomous Trees. Within these trees, polytomies are resolved by randomly coalescing two offspring
34 until binary trees, termed Genetic Resolved Trees, are obtained. These genetic distances are then re-dated

1 using LSD2 (To et al., 2016), assigning dates to the tips by adding the lengths from the tips to the root within
2 the time-scaled trees to a dummy date designated as the root date. Additionally, a temporal constraint for
3 the root is established by setting a range (dummy date - 1 day, dummy date + 1 day), ensuring the root's
4 time is not excessively early. The clock rate used is the same as mentioned above, with a mean of 8×10^{-4}
5 and a standard deviation of 4×10^{-4} substitutions per site per year.

6 Additional 100,000 trees were simulated, and the PhyloDeep methodology was applied to establish the
7 95% CIs.

8 9 **Feature representation and neural network models**

10 We represent time-scaled phylogenetic trees using sampling probability and 98 SSs, as employed in
11 PhyloDeep (Saulnier et al., 2017; Voznica et al., 2022). However, for genetic distance trees, certain concepts
12 like transmission chains (14 SSs) associated with superspreading and lineage through time (LTT) (49 SSs) are
13 not directly applicable. To address this, we designed 62 SSs to capture the distribution of nodes in the
14 phylogenetic tree: 31 SSs for internal nodes (non-leaf nodes within the tree structure, corresponding to
15 transmission events), and 31 SSs for external nodes (leaves of the tree, corresponding to sampling events),
16 by counting the nodes that are n (0-30) mutations away from the tree root. Additionally, we included 10
17 summary statistics related to the size distribution of clusters of identical sequences. These counts capture the
18 number of clusters for each size from 1 to 9, with a combined count for clusters larger than 9, reflecting the
19 underlying transmission dynamics and heterogeneity (Tran-Kiem & Bedford, 2024). Consequently, 90 SS are
20 utilized to characterize the genetic distance tree. While time-scaled trees are rescaled so the average
21 branch length equals 1 prior to representation (Voznica et al., 2022), genetic distance trees do not require
22 this adjustment.

23 Following the PhyloDeep methodology, we implemented our neural network model using Python 3.6, with
24 the Tensorflow 1.5.0, Keras 2.2.4, and scikit-learn 0.19.1 libraries. We partitioned 200,000 simulated
25 phylogenetic trees into 190,000 for training, 9,000 for validation, and 1,000 for testing. The network
26 architecture includes an input layer with either 99 or 90 nodes, followed by four sequential hidden layers
27 arranged in a funnel shape with 64, 32, 16, and 8 neurons, respectively, and an output layer that predicts
28 the four parameters of the BDSS model (R_0 , $1/\gamma$, X_{ss} , and f_{ss}). We experimented with adding or removing
29 hidden layers in the Baseline-Model, which did not improve accuracy. The neurons in the last hidden layer
30 utilize linear activation, whereas the others employ exponential linear (ELU) activation. The model employs
31 the Adam optimization algorithm and uses Mean Absolute Percentage Error (MAPE) as the loss function, with
32 a batch size of 200 and a maximum of 1,000 epochs. Early stopping, with a patience value of 50, was used
33 to prevent overfitting based on MAPE performance on the validation set. A dropout rate of 0.5 was applied
34 in the hidden layers, and variations in dropout rates between 0.3 and 0.7 did not enhance the Baseline-

1 Model's accuracy. The performance of our neural network models is assessed as the mean relative error
 2 (MRE) of the estimator:

$$3 \quad MRE = \frac{1}{n} \sum_{i=1}^n \left(\frac{\text{predicted}_i - \text{target}_i}{\text{target}_i} \right)$$

4 where n is the number of simulated trees used in the test set.

5 To draw a parallel with epidemiological inference, X_{ss} and f_{ss} can be transformed into the dispersion k .
 6 Utilizing the multi-type birth-death model process (Stadler & Bonhoeffer, 2013), it becomes possible to
 7 estimate the probability of an individual infecting “ n ” others over its lifespan, aligning with a geometric
 8 distribution. By synthesizing the probability with the cumulative number of infections, the offspring distribution
 9 was ascertained. The approach outlined in “Estimating R_0 and k from epidemiological data only” section was
 10 employed to derive k from this offspring distribution.

11

12 **Integration of contact tracing data into phylogenetic trees**

13 In our simulations, we utilize time-scaled trees to derive contact tracing data, treating all descendants of
 14 each internal node as a single cluster, with the node's age representing the infection time (Supplementary
 15 Figure S3). Using such contact tracing data, we refine the phylogenetic trees by identifying the most recent
 16 common ancestor (MRCA) for each cluster. We then iterate through children of the MRCA and coalesce all
 17 associated children, encompassing both leaves and children of internal nodes within the cluster. This process
 18 enables us to resolve polytomies in Genetic Polytomous trees, facilitating their transformation back into
 19 Genetic Baseline trees (Supplementary Figure S2).

20 Additionally, by applying the infection times as time constraints on the internal nodes, we can revert Genetic
 21 Baseline trees to their Baseline counterparts using LSD2 (To et al., 2016). We achieve this by setting a
 22 specific time range for the internal nodes, using a margin of (infection time - 1 day, infection time + 1 day).
 23 Narrowing this margin to 0.1 day brings the converted trees even closer to the Baseline trees, thereby
 24 yielding performance on the Baseline-Model that is nearly identical to that obtained when directly using
 25 Baseline trees for testing, as detailed in Supplementary Tables S2 and S3.

26

27 **SARS-CoV-2 dataset in Hong Kong**

28 We used sequences and epidemiological data from the third and fourth waves of SARS-CoV-2 in Hong
 29 Kong, as detailed in our prior study (Gu et al., 2022). These waves were characterized by single introduction
 30 events that sparked local transmissions, and they were notable for their relatively consistent sequence
 31 sampling and comprehensive surveillance data. In this study, we focused on the exponential stages of waves
 32 3 and 4, which spanned from May 13 to August 1, 2020, with 460 sequences and 1,930 local cases, and

1 from September 30 to December 8, 2020, with 243 sequences and 1,577 local cases, respectively. The
2 sampling rates for waves 3 and 4 were 23.8% and 15.4%, respectively. During wave 3, 84.35% (388 out
3 of 460) of sequences were linked to cluster information involving 191 clusters, among which 76 clusters
4 comprised more than one sequence. This indicates that 16.56% (76 out of 459) of the data were supported
5 by contact tracing. In wave 4, 90.53% (220 out of 243) of sequences were associated with 35 clusters, with
6 23 clusters containing multiple sequences, amounting to 9.50% (23 out of 242) contact tracing data
7 availability.

8 For waves 3 and 4, we reconstructed Maximum Likelihood (ML) phylogenies using RAxML-NG (Kozlov et al.,
9 2019) with the GTR+G4+FO substitution model. We maintained consistency with simulated trees in terms of
10 collapsing internal nodes and the random resolution of polytomies. Our findings revealed that the distribution
11 of the number of offspring from collapsed internal nodes falls within the range observed in our simulations
12 (Supplementary Figure S1). Subsequently, these trees were dated using LSD2 (To et al., 2016), following a
13 strict molecular clock assumption of 8×10^{-4} substitutions per site per year (Hadfield et al., 2018; Jolly &
14 Scaria, 2021), and applying time constraints for the root as inferred by (Gu et al., 2022).

15 16 **Estimating R_0 and k from epidemiological data only**

17 We compared the results for R_0 and k estimated using our deep learning models to those estimated from
18 line-list data on SARS-CoV-2 available during the exponential periods of waves 3 and 4 in Hong Kong.
19 Comparable estimates of R_0 were estimated as per methods described in (Wallinga & Lipsitch, 2007) and
20 implemented in the R package R_0 (Obadia et al., 2012) which assumes an SIR model of transmission like
21 BDSS. We used line-listed incidence data of SARS-CoV-2 symptom onset dates and an exponential
22 generation time distribution also like BDSS (mean = 5.7, SD = 1.8 (Hu et al., 2021)) with results listed in
23 **Table 1**. Additional sensitivity analyses were conducted assuming alternative parameterisations of the
24 generation time (mean = 7.27, SD = 3.81 (Chen et al., 2022)), and/or a gamma-distributed generation
25 time are summarized in Supplementary Table S8.

26 Epidemiological estimates of k were generated by constructing empirical offspring distributions from contact
27 tracing data on SARS-CoV-2 available from previous studies in Hong Kong (Adam et al., 2022). These
28 distributions were generated from infector-infectee pairs, where the number of secondary cases is counted
29 for each unique infector and includes chain-terminating infectees as zero. We subsetted the empirical
30 offspring distributions to the same exponential periods for wave 3 and wave 4 as before, given the
31 estimated infection date of each paired case as a deconvolution of the generation time, incubation period,
32 and delay distributions given the onset date or report dates if asymptomatic between infector-infectee pairs.
33 Importantly, offspring counts were not artificially right-censored, meaning the observed count of each infector
34 case was included even if the estimated infection date of paired infectee(s) fell outside the exponential
35 periods of each wave. Following the approach of Lloyd-Smith et al (Lloyd-Smith et al., 2005), k is estimated

1 directly from the finalised offspring distributions by maximum likelihood estimation, assuming a negative
2 binomial model jointly parameterised by the mean and dispersion parameter k , with 95% intervals
3 generated by non-parametric bootstrap estimation sampling 1000 replicates with replacement.

4 5 **Parameter inference comparison with BEAST2**

6 We assessed the predictive performance of the Dated Baseline-Model and Baseline-Model against the well-
7 established Bayesian structured birth-death model, implemented via the *bdmm* package (Scire et al., 2020)
8 in BEAST2 (Bouckaert et al., 2014) (version 2.6.2). We applied the same priors as used in PhyloDeep
9 (Voznica et al., 2022), maintaining the equality $\beta_{SS}/\beta_{NS} = \beta_{SN}/\beta_{NN}$ and fixing the sampling proportion
10 and tree topology during parameter estimation. Markov Chain Monte Carlo (MCMC) analysis was run for
11 10 million steps, sampling every 1,000 steps with a 10% as burn-in, and Effective Sample Size (ESS) values
12 were assessed using Tracer (Rambaut et al., 2018). The analysis was conducted on 100 simulated Genetic
13 Polytomous Trees incorporating varying levels of contact tracing data (0%, 50%, and 100%) to facilitate
14 transforming the input trees back into Baseline and Dated Baseline Trees, the latter using only cluster
15 information. Additionally, we conducted the BEAST2 analysis on the Hong Kong datasets, which produced
16 different estimations (Supplementary Table S9). However, the poor performance in our simulation analysis
17 without contact tracing data, or when only incorporating cluster information, along with the limited cluster
18 data available in the Hong Kong datasets, was insufficient to meaningfully improve predictions.

19
20 **Acknowledgments:** We acknowledge the technical support provided by colleagues from the Centre for
21 PanorOmic Sciences of the University of Hong Kong. We also acknowledge the Centre for Health Protection
22 of the Department of Health for providing epidemiological data for the study. The computations were
23 performed using research computing facilities offered by Information Technology Services, the University of
24 Hong Kong. The funding bodies had no role in the design of the study and collection, analysis, and
25 interpretation of data and writing of the manuscript. The work described in this paper was substantially
26 supported by a fellowship award from the Research Grants Council of the Hong Kong Special Administrative
27 Region, China (Project No. HKU PDFS2425-7S01).

28 29 **Funding:**

30 National Institutes of Health contract number 75N93021C00016 (VD)

31 Research Grants Council of the Hong Kong SAR, China (Project No. [T11-705/21-N]) (VD)

32 The Collaborative Research Scheme (Project No. C7123-20G) of the Research Grants Council of the Hong
33 Kong Special Administrative Region, China (BC, DA)

1 Health and Medical Research Fund Seed Grant Scheme (Project No. 22211192) of the Hong Kong SAR
2 (DA)
3 HKU-Pasteur Research Pole Fellowship 2023 (S-AC23005-01) (RX)
4 Research Grants Council of the Hong Kong SAR, China (Project No. [HKU PDFS2425-7S01]) (RX)
5 PaRis AI Research InstitutE (PRAIRIE; ANR-19-P3IA-0001) (OG)
6

7 **Competing interests:** Authors declare that they have no competing interests.

8
9 **Data and materials availability:** All anonymized data, code, and analysis files are available in the
10 GitHub repository (<https://github.com/vjlab/dl-phylogenetics-ct>).

11 Reference

- 12 Adam, D., Gostic, K., Tsang, T., Wu, P., Lim, W. W., Yeung, A., . . . Chen, D. (2022). Time-varying
13 transmission heterogeneity of SARS and COVID-19 in Hong Kong. *Research Square*.
- 14 Adam, D. C., Wu, P., Wong, J. Y., Lau, E. H. Y., Tsang, T. K., Cauchemez, S., . . . Cowling, B. J.
15 (2020). Clustering and superspreading potential of SARS-CoV-2 infections in Hong Kong.
16 *Nat Med*, 26(11), 1714-1719. <https://doi.org/10.1038/s41591-020-1092-0>
- 17 Bouckaert, R., Heled, J., Kuhnert, D., Vaughan, T., Wu, C. H., Xie, D., . . . Drummond, A. J. (2014).
18 BEAST 2: a software platform for Bayesian evolutionary analysis. *PLoS Comput Biol*,
19 10(4), e1003537. <https://doi.org/10.1371/journal.pcbi.1003537>
- 20 Chen, D., Lau, Y. C., Xu, X. K., Wang, L., Du, Z., Tsang, T. K., . . . Ali, S. T. (2022). Inferring time-
21 varying generation time, serial interval, and incubation period distributions for COVID-
22 19. *Nat Commun*, 13(1), 7727. <https://doi.org/10.1038/s41467-022-35496-8>
- 23 Drummond, A. J., Rambaut, A., Shapiro, B., & Pybus, O. G. (2005). Bayesian coalescent inference
24 of past population dynamics from molecular sequences. *Mol Biol Evol*, 22(5), 1185-1192.
25 <https://doi.org/10.1093/molbev/msi103>
- 26 Du, Z., Wang, C., Liu, C., Bai, Y., Pei, S., Adam, D. C., . . . Cowling, B. J. (2022). Systematic review
27 and meta-analyses of superspreading of SARS-CoV-2 infections. *Transbound Emerg Dis*.
28 <https://doi.org/10.1111/tbed.14655>
- 29 Eden, J. S., Sikazwe, C., Xie, R., Deng, Y. M., Sullivan, S. G., Michie, A., . . . Australian, R. S. V. s. g.
30 (2022). Off-season RSV epidemics in Australia after easing of COVID-19 restrictions. *Nat*
31 *Commun*, 13(1), 2884. <https://doi.org/10.1038/s41467-022-30485-3>
- 32 Featherstone, L. A., Duchene, S., & Vaughan, T. G. (2023). Decoding the Fundamental Drivers of
33 Phylodynamic Inference. *Mol Biol Evol*, 40(6). <https://doi.org/10.1093/molbev/msad132>
- 34 Geidelberg, L., Boyd, O., Jorgensen, D., Siveroni, I., Nascimento, F. F., Johnson, R., . . . Nie, Q.
35 (2021). Genomic epidemiology of a densely sampled COVID-19 outbreak in China. *Virus*
36 *Evol*, 7(1), veaa102. <https://doi.org/10.1093/ve/veaa102>
- 37 Gu, H., Xie, R., Adam, D. C., Tsui, J. L., Chu, D. K., Chang, L. D. J., . . . Poon, L. L. M. (2022).
38 Genomic epidemiology of SARS-CoV-2 under an elimination strategy in Hong Kong. *Nat*
39 *Commun*, 13(1), 736. <https://doi.org/10.1038/s41467-022-28420-7>

- 1 Guindon, S., Dufayard, J. F., Lefort, V., Anisimova, M., Hordijk, W., & Gascuel, O. (2010). New
2 algorithms and methods to estimate maximum-likelihood phylogenies: assessing the
3 performance of PhyML 3.0. *Syst Biol*, 59(3), 307-321.
4 <https://doi.org/10.1093/sysbio/syq010>
- 5 Guo, Z., Zhao, S., Lee, S. S., Mok, C. K. P., Wong, N. S., Wang, J., . . . Yeoh, E. K. (2022).
6 Superspreading potential of COVID-19 outbreak seeded by Omicron variants of SARS-
7 CoV-2 in Hong Kong. *J Travel Med*. <https://doi.org/10.1093/jtm/taac049>
- 8 Hadfield, J., Megill, C., Bell, S. M., Huddleston, J., Potter, B., Callender, C., . . . Neher, R. A.
9 (2018). Nextstrain: real-time tracking of pathogen evolution. *Bioinformatics*, 34(23),
10 4121-4123. <https://doi.org/10.1093/bioinformatics/bty407>
- 11 Hohna, S., & Drummond, A. J. (2012). Guided tree topology proposals for Bayesian phylogenetic
12 inference. *Syst Biol*, 61(1), 1-11. <https://doi.org/10.1093/sysbio/syr074>
- 13 Hu, S., Wang, W., Wang, Y., Litvinova, M., Luo, K., Ren, L., . . . Yu, H. (2021). Infectivity,
14 susceptibility, and risk factors associated with SARS-CoV-2 transmission under intensive
15 contact tracing in Hunan, China. *Nat Commun*, 12(1), 1533.
16 <https://doi.org/10.1038/s41467-021-21710-6>
- 17 Jolly, B., & Scaria, V. (2021). Computational Analysis and Phylogenetic Clustering of SARS-CoV-2
18 Genomes. *Bio Protoc*, 11(8), e3999. <https://doi.org/10.21769/BioProtoc.3999>
- 19 Kozlov, A. M., Darriba, D., Flouri, T., Morel, B., & Stamatakis, A. (2019). RAxML-NG: a fast,
20 scalable and user-friendly tool for maximum likelihood phylogenetic inference.
21 *Bioinformatics*, 35(21), 4453-4455. <https://doi.org/10.1093/bioinformatics/btz305>
- 22 Kuhnert, D., Stadler, T., Vaughan, T. G., & Drummond, A. J. (2016). Phylodynamics with
23 Migration: A Computational Framework to Quantify Population Structure from Genomic
24 Data. *Mol Biol Evol*, 33(8), 2102-2116. <https://doi.org/10.1093/molbev/msw064>
- 25 Lambert, S., Voznica, J., & Morlon, H. (2023). Deep Learning from Phylogenies for Diversification
26 Analyses. *Syst Biol*. <https://doi.org/10.1093/sysbio/syad044>
- 27 Leventhal, G. E., Gunthard, H. F., Bonhoeffer, S., & Stadler, T. (2014). Using an epidemiological
28 model for phylogenetic inference reveals density dependence in HIV transmission. *Mol*
29 *Biol Evol*, 31(1), 6-17. <https://doi.org/10.1093/molbev/mst172>
- 30 Lewis, P. O., Holder, M. T., & Holsinger, K. E. (2005). Polytomies and Bayesian phylogenetic
31 inference. *Syst Biol*, 54(2), 241-253. <https://doi.org/10.1080/10635150590924208>
- 32 Lloyd-Smith, J. O., Schreiber, S. J., Kopp, P. E., & Getz, W. M. (2005). Superspreading and the
33 effect of individual variation on disease emergence. *Nature*, 438(7066), 355-359.
34 <https://doi.org/10.1038/nature04153>
- 35 Miller, D., Martin, M. A., Harel, N., Tirosh, O., Kustin, T., Meir, M., . . . Stern, A. (2020). Full
36 genome viral sequences inform patterns of SARS-CoV-2 spread into and within Israel.
37 *Nat Commun*, 11(1), 5518. <https://doi.org/10.1038/s41467-020-19248-0>
- 38 Morel, B., Barbera, P., Czech, L., Bettisworth, B., Hubner, L., Lutteropp, S., . . . Stamatakis, A.
39 (2021). Phylogenetic Analysis of SARS-CoV-2 Data Is Difficult. *Mol Biol Evol*, 38(5), 1777-
40 1791. <https://doi.org/10.1093/molbev/msaa314>
- 41 Moreno, G. K., Braun, K. M., Riemersma, K. K., Martin, M. A., Halfmann, P. J., Crooks, C. M., . . .
42 Friedrich, T. C. (2020). Revealing fine-scale spatiotemporal differences in SARS-CoV-2
43 introduction and spread. *Nat Commun*, 11(1), 5558. [https://doi.org/10.1038/s41467-
44 020-19346-z](https://doi.org/10.1038/s41467-020-19346-z)

- 1 Nguyen, L. T., Schmidt, H. A., von Haeseler, A., & Minh, B. Q. (2015). IQ-TREE: a fast and
2 effective stochastic algorithm for estimating maximum-likelihood phylogenies. *Mol Biol*
3 *Evol*, 32(1), 268-274. <https://doi.org/10.1093/molbev/msu300>
- 4 Obadia, T., Haneef, R., & Boelle, P. Y. (2012). The R0 package: a toolbox to estimate
5 reproduction numbers for epidemic outbreaks. *BMC Med Inform Decis Mak*, 12, 147.
6 <https://doi.org/10.1186/1472-6947-12-147>
- 7 Paredes, M. I., Ahmed, N., Figgins, M., Colizza, V., Lemey, P., McCrone, J. T., . . . Bedford, T.
8 (2024). Underdetected dispersal and extensive local transmission drove the 2022 mpox
9 epidemic. *Cell*. <https://doi.org/10.1016/j.cell.2024.02.003>
- 10 Pekar, J. E., Magee, A., Parker, E., Moshiri, N., Izhikevich, K., Havens, J. L., . . . Wertheim, J. O.
11 (2022). The molecular epidemiology of multiple zoonotic origins of SARS-CoV-2. *Science*,
12 377(6609), 960-966. <https://doi.org/10.1126/science.abp8337>
- 13 Price, M. N., Dehal, P. S., & Arkin, A. P. (2010). FastTree 2--approximately maximum-likelihood
14 trees for large alignments. *PLoS One*, 5(3), e9490.
15 <https://doi.org/10.1371/journal.pone.0009490>
- 16 Ragonnet-Cronin, M., Boyd, O., Geidelberg, L., Jorgensen, D., Nascimento, F. F., Siveroni, I., . . .
17 Volz, E. (2021). Genetic evidence for the association between COVID-19 epidemic
18 severity and timing of non-pharmaceutical interventions. *Nat Commun*, 12(1), 2188.
19 <https://doi.org/10.1038/s41467-021-22366-y>
- 20 Rambaut, A., Drummond, A. J., Xie, D., Baele, G., & Suchard, M. A. (2018). Posterior
21 Summarization in Bayesian Phylogenetics Using Tracer 1.7. *Syst Biol*, 67(5), 901-904.
22 <https://doi.org/10.1093/sysbio/syy032>
- 23 Sagulenko, P., Puller, V., & Neher, R. A. (2018). TreeTime: Maximum-likelihood phylodynamic
24 analysis. *Virus Evol*, 4(1), vex042. <https://doi.org/10.1093/ve/vex042>
- 25 Saulnier, E., Gascuel, O., & Alizon, S. (2017). Inferring epidemiological parameters from
26 phylogenies using regression-ABC: A comparative study. *PLoS Comput Biol*, 13(3),
27 e1005416. <https://doi.org/10.1371/journal.pcbi.1005416>
- 28 Scire, J., Barido-Sottani, J., Kühnert, D., Vaughan, T. G., & Stadler, T. (2020). Improved multi-
29 type birth-death phylodynamic inference in BEAST 2. *BioRxiv*, 2020.2001. 2006.895532.
- 30 Stadler, T., & Bonhoeffer, S. (2013). Uncovering epidemiological dynamics in heterogeneous
31 host populations using phylogenetic methods. *Philos Trans R Soc Lond B Biol Sci*,
32 368(1614), 20120198. <https://doi.org/10.1098/rstb.2012.0198>
- 33 Stadler, T., Kouyos, R., von Wyl, V., Yerly, S., Boni, J., Burgisser, P., . . . Swiss, H. I. V. C. S. (2012).
34 Estimating the basic reproductive number from viral sequence data. *Mol Biol Evol*, 29(1),
35 347-357. <https://doi.org/10.1093/molbev/msr217>
- 36 Stadler, T., Kuhnert, D., Bonhoeffer, S., & Drummond, A. J. (2013). Birth-death skyline plot
37 reveals temporal changes of epidemic spread in HIV and hepatitis C virus (HCV). *Proc*
38 *Natl Acad Sci U S A*, 110(1), 228-233. <https://doi.org/10.1073/pnas.1207965110>
- 39 Thompson, A., Liebeskind, B., Scully, E. J., & Landis, M. (2024). Deep learning and likelihood
40 approaches for viral phylogeography converge on the same answers whether the
41 inference model is right or wrong. *Syst Biol*. <https://doi.org/10.1093/sysbio/syad074>
- 42 To, T. H., Jung, M., Lycett, S., & Gascuel, O. (2016). Fast Dating Using Least-Squares Criteria and
43 Algorithms. *Syst Biol*, 65(1), 82-97. <https://doi.org/10.1093/sysbio/syv068>

- 1 Tran-Kiem, C., & Bedford, T. (2024). Estimating the reproduction number and transmission
 2 heterogeneity from the size distribution of clusters of identical pathogen sequences.
 3 *Proc Natl Acad Sci U S A*, *121*(15), e2305299121.
 4 <https://doi.org/10.1073/pnas.2305299121>
- 5 Vaughan, T. G., Scire, J., Nadeau, S. A., & Stadler, T. (2024). Estimates of early outbreak-specific
 6 SARS-CoV-2 epidemiological parameters from genomic data. *Proc Natl Acad Sci U S A*,
 7 *121*(2), e2308125121. <https://doi.org/10.1073/pnas.2308125121>
- 8 Volz, E. M., Kosakovsky Pond, S. L., Ward, M. J., Leigh Brown, A. J., & Frost, S. D. (2009).
 9 Phylodynamics of infectious disease epidemics. *Genetics*, *183*(4), 1421-1430.
 10 <https://doi.org/10.1534/genetics.109.106021>
- 11 Volz, E. M., & Siveroni, I. (2018). Bayesian phylodynamic inference with complex models. *PLoS*
 12 *Comput Biol*, *14*(11), e1006546. <https://doi.org/10.1371/journal.pcbi.1006546>
- 13 Voznica, J., Zhukova, A., Boskova, V., Saulnier, E., Lemoine, F., Moslonka-Lefebvre, M., &
 14 Gascuel, O. (2022). Deep learning from phylogenies to uncover the epidemiological
 15 dynamics of outbreaks. *Nat Commun*, *13*(1), 3896. [https://doi.org/10.1038/s41467-022-](https://doi.org/10.1038/s41467-022-31511-0)
 16 [31511-0](https://doi.org/10.1038/s41467-022-31511-0)
- 17 Wallinga, J., & Lipsitch, M. (2007). How generation intervals shape the relationship between
 18 growth rates and reproductive numbers. *Proc Biol Sci*, *274*(1609), 599-604.
 19 <https://doi.org/10.1098/rspb.2006.3754>
- 20 Wells, C. R., Townsend, J. P., Pandey, A., Moghadas, S. M., Krieger, G., Singer, B., . . . Galvani, A.
 21 P. (2021). Optimal COVID-19 quarantine and testing strategies. *Nat Commun*, *12*(1), 356.
 22 <https://doi.org/10.1038/s41467-020-20742-8>
- 23 Worobey, M., Han, G. Z., & Rambaut, A. (2014). Genesis and pathogenesis of the 1918
 24 pandemic H1N1 influenza A virus. *Proc Natl Acad Sci U S A*, *111*(22), 8107-8112.
 25 <https://doi.org/10.1073/pnas.1324197111>
- 26 Zhukova, A., & Gascuel, O. (2024). Accounting for partner notification in epidemiological birth-
 27 death-models. *medRxiv*, 2024.2009.2009.24313296.

30 Figure Legend

31 **Fig. 1. An overview of training neural network models based on simulated phylogenetic trees.** The BDSS model
 32 categorizes individuals as superspreaders (S) or normal spreaders (N), extending the traditional Birth-Death model
 33 parameters R_0 (basic reproductive number) and $1/\gamma$ (infectious period). f_{ss} indicates the fraction of superspreaders in
 34 the population, while X_{ss} represents the ratio of the transmission rate of superspreaders to that of normal spreaders.
 35 Seven types of trees, Baseline, Dated Baseline, Dated Resolved, Dated Polytomous, Genetic Baseline, Genetic Resolved,
 36 Genetic Polytomous, are detailed in Fig. 2.

37 **Fig. 2. Examples of seven types of phylogenetic trees used in simulations.** Internal nodes are marked as black dots,
 38 while tips are denoted by numerical labels. Among these, four trees represent poorly resolved, realistic phylogenetic
 39 structures that can be derived from sequence data and are highlighted with a grey background. To effectively highlight
 40 the differences between poorly resolved trees, which can be constructed from sequence data, and fully resolved
 41 idealistic trees, which cannot, tips have been color-coded into three distinct clusters. Each type of simulated tree used in
 42 this study has tip counts ranging from 200 to 1000 (Supplementary Table S1).

1 **Fig. 3. Performance comparison of models.** A) Performance comparison of models trained on seven types of
 2 phylogenetic trees. Each bar depicts the relative error observed when testing trees of the same type as those used in
 3 training. The red marked lines denote the median relative error when testing the Baseline-Model and Dated Baseline-
 4 Model with Dated Resolved trees, as well as the Genetic Baseline-Model with Genetic Resolved trees. Models trained
 5 using poorly resolved phylogenetic trees (i.e., Dated Resolved, Dated Polytomous, Genetic Resolved and Genetic
 6 Polytomous) are highlighted in bold. B) Performance comparison of models tested using poorly resolved phylogenetic
 7 trees. "Baseline-Poor" represents the evaluation of the Baseline-Model tested using Dated Resolved Trees. "Dated
 8 Baseline-Poor" indicates the assessment of the Dated Baseline-Model with Dated Resolved Trees, while "Genetic
 9 Baseline-Poor" reflects the performance of the Genetic Baseline-Model when testing with Genetic Resolved trees.

10 **Fig. 4. Performance comparison by incorporating varying levels of contact tracing data based on Baseline-Model,**
 11 **Dated Baseline-Model and BEAST2.** A) Comparison between the Baseline-Model and Dated Baseline-Model with
 12 varying levels of contact tracing data based on 1,000 simulated trees. The models are represented by grey (Baseline-
 13 Model) and red (Dated Baseline-Model) bars, with the color intensity within each bar signaling the degree of contact
 14 tracing data integrated into the input trees. Darker shades denote a higher percentage of data incorporation. The term
 15 "Baseline_50" refers to the performance of the Baseline-Model with Genetic Polytomous trees refined using 50%
 16 contact tracing data, encompassing cluster information and infection times. "Dated Baseline_50" indicates the
 17 performance of the Dated Baseline-Model with Genetic Polytomous trees refined using 50% contact tracing data,
 18 including only cluster information. It's notable that the input trees are refined by infection time, with a 1-day time
 19 constraint margin using LSD2 (To et al., 2016), and an additional refinement with a stricter margin of 0.1 day, as shown
 20 in Supplementary Table S3. B) Comparison between the Baseline-Model and BEAST2 (blue bar) with varying levels of
 21 contact tracing data (cluster information and infection times) based on 100 simulated trees. C) Comparison between
 22 Dated Baseline-Model and BEAST2 (blue bar) with varying levels of contact tracing data (cluster information only)
 23 based on 100 simulated trees. "BEAST2_50" indicates the performance of BEAST2 with Genetic Polytomous trees
 24 refined using 50% contact tracing data, incorporating both cluster information and infection times in panel B, and only
 25 cluster information in panel C.

26
 27
 28
 29

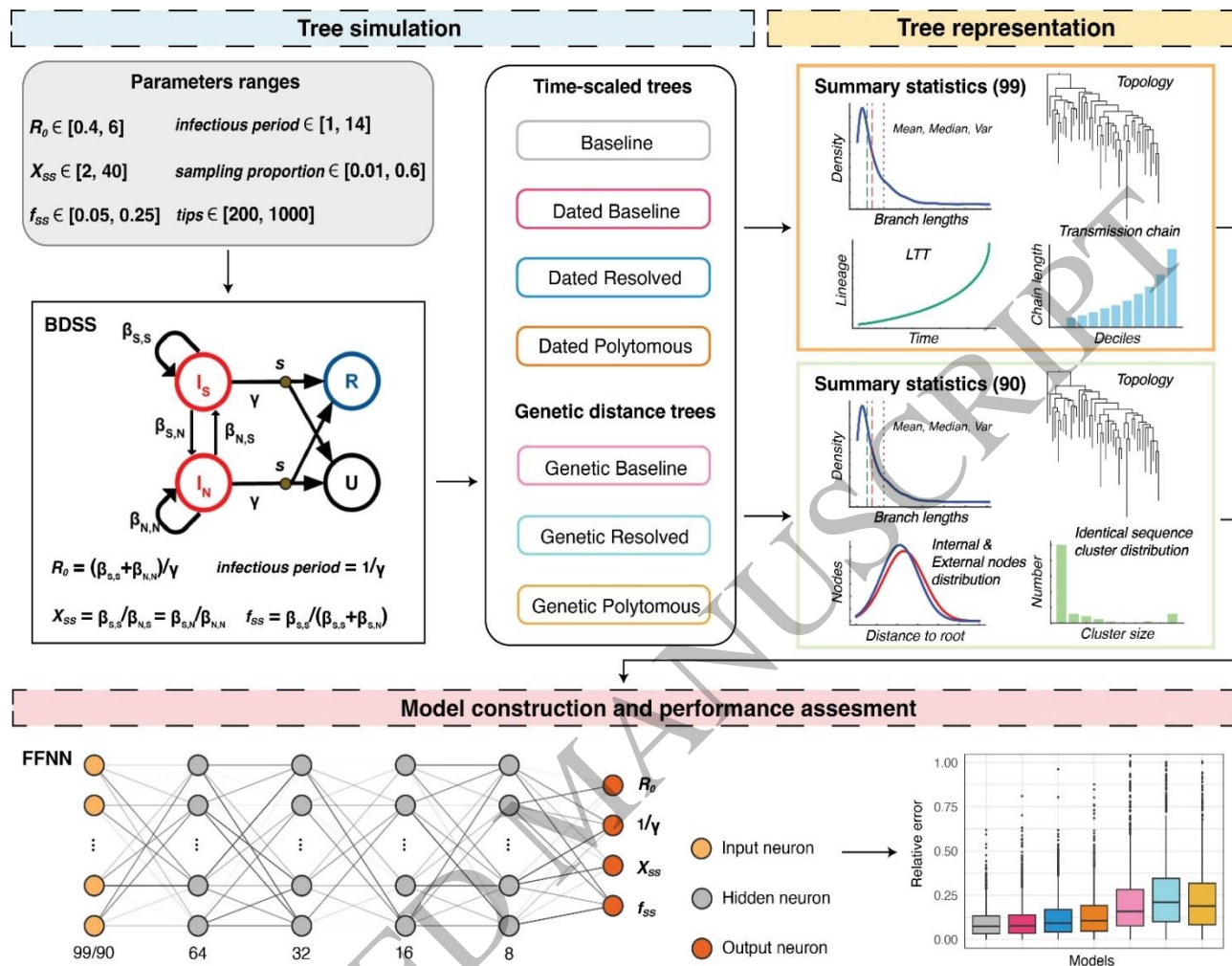


Figure 1
184x139 mm (x DPI)

1

2

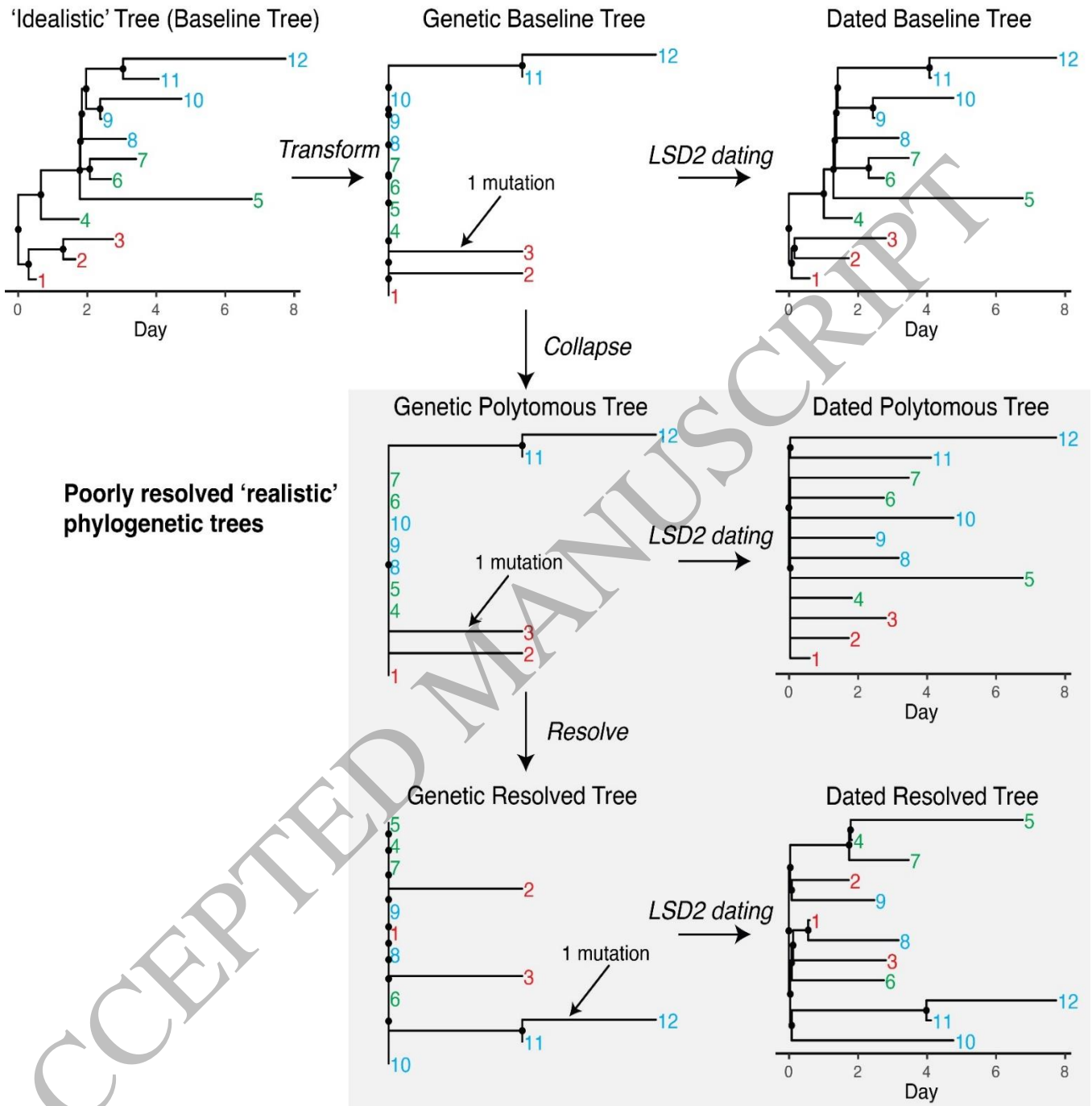


Figure 2
203x187 mm (x DPI)

1

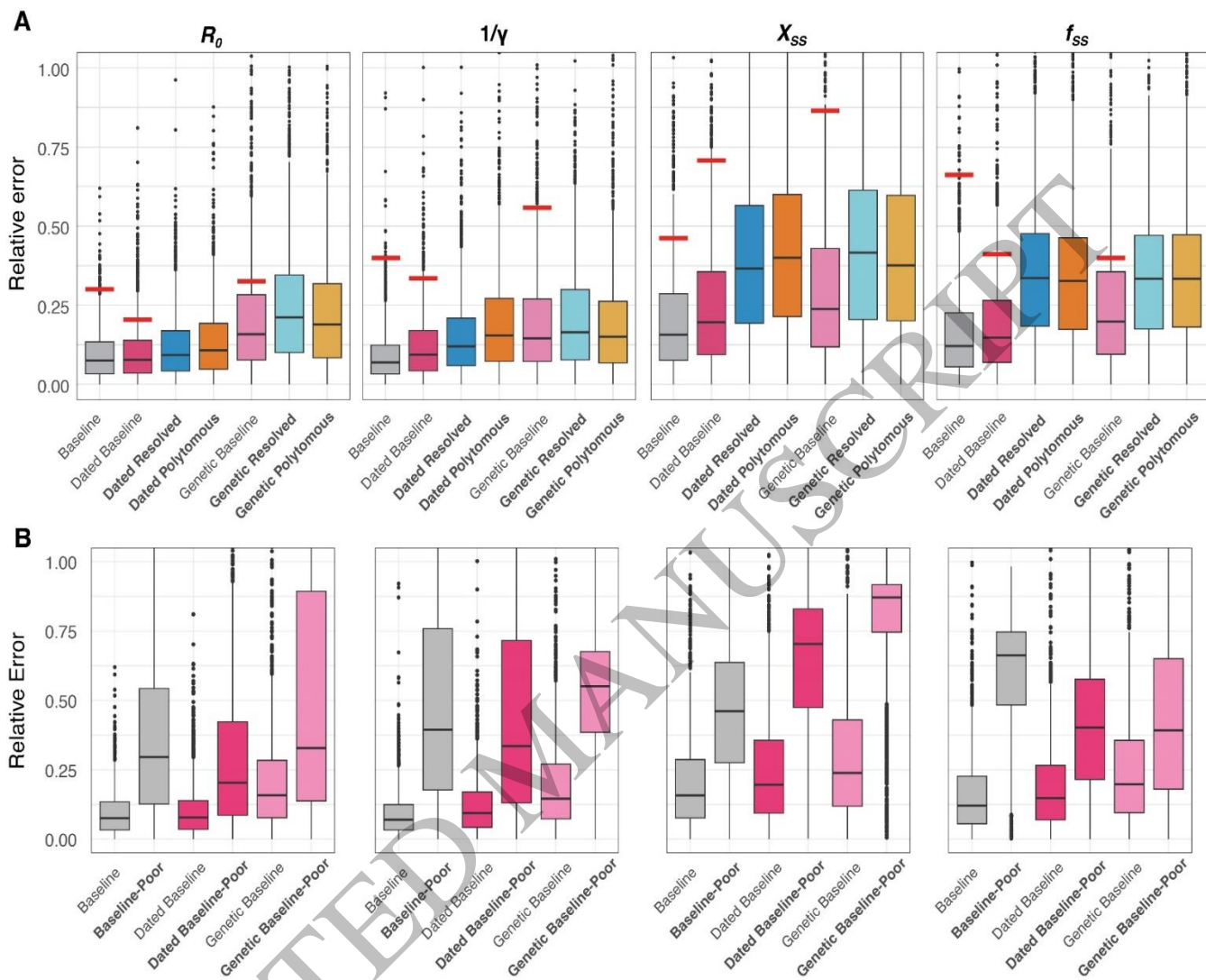


Figure 3
206x153 mm (x DPI)

1

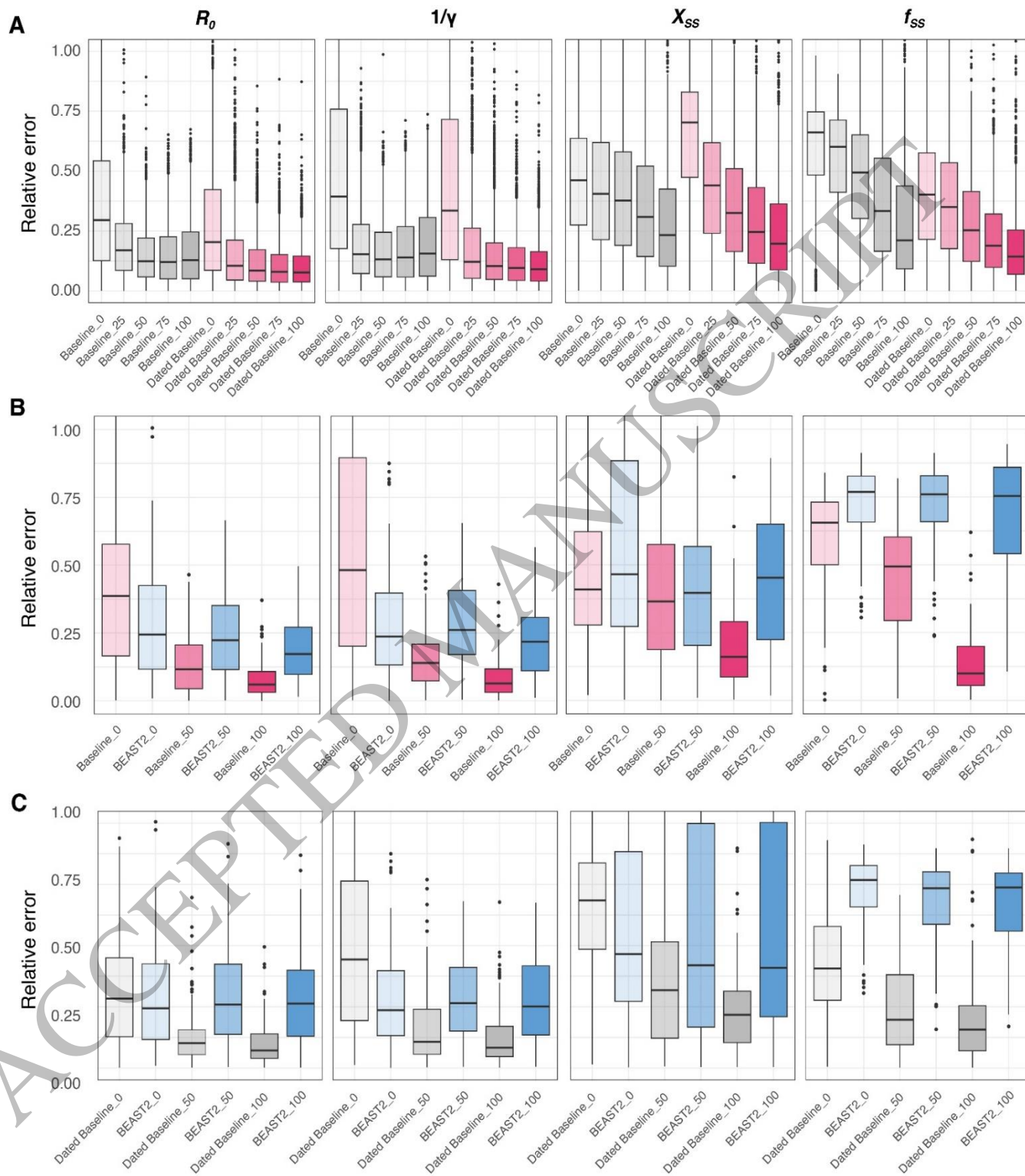


Figure 4
188x209 mm (x DPI)