

 Open access • Posted Content • DOI:10.1101/2020.02.10.942367

## Polynomial Phylogenetic Analysis of Tree Shapes — Source link

[Pengyu Liu](#), [Priscila Biller](#), [Matthew Gould](#), [Caroline Colijn](#)

**Institutions:** [Simon Fraser University](#)

**Published on:** 29 Sep 2020 - [bioRxiv](#) (Cold Spring Harbor Laboratory)

**Topics:** [Tree \(data structure\)](#) and [Phylogenetic tree](#)

Related papers:

- [A Metric on Phylogenetic Tree Shapes.](#)
- [Network science inspires novel tree shape statistics](#)
- [Discrete coalescent trees.](#)
- [Tree-based networks: characterisations, metrics, and support trees](#)
- [Optimization over a class of tree shape statistics](#)

Share this paper:    

View more about this paper here: <https://typeset.io/papers/polynomial-phylogenetic-analysis-of-tree-shapes-13q5izgw2b>

# Polynomial Phylogenetic Analysis of Tree Shapes

PENGYU LIU\*, PRISCILA BILLER, MATTHEW GOULD, AND CAROLINE COLIJN,  
*Department of Mathematics, Simon Fraser University, Burnaby, V5A 1S6, Canada*

*\*E-mail: pengyu\_liu@sfu.ca*

## ABSTRACT

1 Phylogenetic trees are a central tool in evolutionary biology. They demonstrate  
2 evolutionary patterns among species, genes, and with modern sequencing technologies,  
3 patterns of ancestry among sets of individuals. Phylogenetic trees usually consist of tree  
4 shapes, branch lengths and partial labels. Comparing tree shapes is a challenging aspect of  
5 comparing phylogenetic trees as there are few tools to describe tree shapes in a  
6 quantitative, accurate, comprehensive and easy-to-interpret way. Current methods to  
7 compare tree shapes are often based on scalar indices reflecting tree imbalance, and on  
8 frequencies of small subtrees. In this paper, we present tree comparisons and applications  
9 based on a polynomial that fully characterizes trees. Polynomials are important tools to  
10 describe discrete structures and have been used to study various objects including graphs  
11 and knots. There are also polynomials that describe rooted trees. We use tree-defining  
12 polynomials to compare tree shapes randomly generated by simulations and tree shapes  
13 reconstructed from data. Moreover, we show that the comparisons can be used to estimate  
14 parameters and to select the best-fit model that generates specific tree shapes.

15 *Key words:* Phylogenetics, Polynomials, Tree Shapes, Tree Metrics

16 A tree is a natural data structure that represents hierarchical relations between  
17 objects. In phylogenetics, a tree structure usually includes its tree shape, that is, the  
18 unlabeled underlying graph, as well as branch lengths reflecting either evolutionary  
19 distance or time. Estimating the branch lengths can be a challenge for tree reconstruction  
20 methods, with Bayesian and maximum likelihood methods yielding inconsistent results  
21 (Brown, 2010), high demands on memory and processor time (Binet, 2016), and/or lack of  
22 strong support for a molecular clock (in the case of timed trees). As a consequence, the  
23 inferred phylogenetic trees may have a consistent tree shape but differing root heights and  
24 branch lengths.

25 The shapes of phylogenetic trees can carry information about macroevolutionary  
26 processes, as well as reflecting the data used and the choice of the evolutionary model  
27 (Kirkpatrick, 1993; Purvis, 2011; Aldous, 1996). The ecological fitness and the presence of  
28 selection can also affect the shapes of trees (Dayarian, 2014; Maia, 2004). In the study of  
29 infectious diseases, where the shapes of phylogenetic trees of pathogens reveal diversity  
30 patterns that represent a combination of unfixated neutral variation, variation under  
31 selection, demographic processes and ecological interactions, it is not clear how informative  
32 the tree shapes are of the underlying evolutionary and epidemiological processes. However,  
33 effort is being made to explore this question, with the main focus often on the frequency of  
34 cherries and tree imbalance (Grenfell, 2004; Lambert, 2013; Plazzotta, 2016; Volz, 2013).

35 One of the main topics of inquiry in phylogenetic tree shapes has been asymmetry,  
36 since a key observation was made that the shapes of phylogenetic trees reconstructed from  
37 data are more asymmetric than tree shapes simulated by simple models (Aldous, 1996).  
38 Various ways to measure the asymmetry were developed (Aldous, 1996; Colless, 1982;  
39 Fusco, 1995; Sackin, 1972; Stich, 2009) and it was shown that these asymmetric measures  
40 can distinguish random trees generated by different models (Agapow, 2002; Kirkpatrick,  
41 1993; Matsen, 2006). At the same time, mathematical models that produce imbalanced  
42 trees were developed (Aldous, 2001; Blum, 2006). As statistical tools, the distributions of

43 tree shapes under simple models can be used to test evolutionary hypotheses (Blum, 2006;  
44 Mooers, 1997; Wu, 2016). In (Manceau, 2015), and mathematical models can be developed  
45 to match the macroevolutionary patterns observed in the phylogenetic trees reconstructed  
46 from data.

47 As the cost of DNA sequencing is decreasing, more genomic data are being collected  
48 and becoming available. More organisms are being sequenced progressively at the  
49 whole-genome scale (Bedford, 2015; Chewapreecha, 2014; Colijn, 2018) and the evolution  
50 of certain pathogens is being tracked in real time (Hadfield, 2018). As a consequence, both  
51 the number and the size of trees reconstructed from data are increasing. Accordingly, a  
52 major challenge in tree shape analysis is that there are few tools to describe and compare  
53 trees in a quantitative, accurate, comprehensive and easy-to-interpret way, especially for  
54 large trees. Scalar indices describing asymmetry or the frequency of subtrees have a  
55 limitation in that many different tree shapes may have the same index. A labelled tree is a  
56 tree shape whose vertices have unique labels. An alternative approach to comparing tree  
57 shapes is using metrics defined for labelled trees, for example, the well known  
58 Robinson-Foulds metric (Robinson, 1981), Billera-Holmes-Vogtmann metric (Billera, 2001)  
59 and Kendall-Colijn metric (Kendall, 2016), among others. These metrics depend on the  
60 labels of the vertices, that is, two labelled trees with the same tree shape but the labels  
61 re-arranged are not identical and the distances between them can be very large. Recently,  
62 metrics defined for rooted unlabelled trees or rooted tree shapes have also been introduced  
63 (Colijn, 2018), making use of integer labels assigned to tree shapes. However, these metrics  
64 have several limitations, including the challenge of interpreting the integer labels, the  
65 treatment of non-binary trees, and the metrics' performance in distinguishing trees from  
66 different processes or datasets.

67 Graph polynomials and knot polynomials are important tools in the mathematical  
68 study of discrete structures, and can be used to describe the structures in interpretable  
69 ways. For example, the Tutte polynomial (Tutte, 1954) is a renowned polynomial for

70 graphs and the Jones polynomial (Jones, 1985) is one of the most important tools to study  
71 knots. In (Liu, 2021), a method to assign a unique polynomial to each tree shape is  
72 introduced. These polynomials provide a new way to describe tree shapes quantitatively  
73 and comprehensively. The coefficients of the polynomial of a tree can be considered as a  
74 generalization of the clade size distribution of the tree. In addition, the set of coefficients of  
75 a tree polynomial can be treated as a vector, and vectors are natural objects on which to  
76 define metrics. In this paper, we introduce the polynomial representations for tree shapes  
77 and we define and examine a metric based on the trees' unique polynomials. We show that  
78 the polynomial representations for tree shapes have perfect resolution and reasonably low  
79 computation time, and the polynomial metric has a performs well at clustering trees,  
80 compared to other high-resolution metrics. We also show that the polynomials can be used  
81 for parameter estimation, and for choosing the best-fit model to generate a tree shape.

## 82 MATERIALS AND METHODS

### 83 *Tree Polynomials*

84 In this paper, a tree shape or simply a tree represents an unlabeled tree, that is a  
85 graph with no cycles, without information about branch lengths or labels unless otherwise  
86 stated. We define the bivariate polynomial  $P(T, x, y)$  for a rooted unlabeled tree  $T$  in the  
87 following way. If  $T$  is the trivial tree with a single vertex, then  $P(T, x, y) = x$ . Otherwise  $T$   
88 has  $k$  branches at its root and each branch leads to a subtree of  $T$ . Let  $T_1, T_2, \dots, T_k$  be  
89 the  $k$  rooted subtrees whose roots are adjacent to the root of  $T$ . We define the polynomial  
90 for  $T$  by  $P(T, x, y) = y + \prod_{i=1}^k P(T_i, x, y)$ . If all of the subtrees are the trivial tree, then the  
91 polynomial is defined and we have a rooted  $k$ -star whose polynomial is  $P(T, x, y) = x^k + y$ .  
92 If there exists a non-trivial subtree  $T_i$ , then we apply the definition to compute  $P(T_i, x, y)$ .  
93 The polynomial  $P(T, x, y)$  can be computed by recursively applying the definition until we  
94 reach all tips of  $T$ . As another example, the polynomial for the three-tip rooted binary tree



*Tree metrics*

114

115 In this paper, we use three tree metrics or distances. The first is a tree metric based  
116 on the Laplacian spectrum. The metric is the Jensen-Shannon distance over the spectrum  
117 densities introduced in (Lewitus, 2016). We call it Lewitus-Morlon metric. The second  
118 metric is based on the subtree size distribution. The subtree size distribution of a tree is  
119 defined as a vector whose  $n$ -th entry is the number of  $n$ -tip subtrees in the tree. The  
120 metric is defined using the Manhattan distance over the subtree size distribution vectors.  
121 We name it the “subtree-Manhatttan metric”. The third metric is based on the  
122 polynomial. Let  $T_1, T_2$  be two trees and  $C(T_1) = (c_1^{(a,b)})$ ,  $C(T_2) = (c_2^{(a,b)})$  be the coefficient  
123 matrices of the polynomials  $P(T_1, x, y)$ ,  $P(T_2, x, y)$ . We define a function

$$124 \quad \mu(c_1, c_2) = \begin{cases} |c_1 - c_2| / (c_1 + c_2) & \text{if } c_1 \neq 0 \text{ or } c_2 \neq 0 \\ 0 & \text{if } c_1 = 0 \text{ and } c_2 = 0 \end{cases}$$

125 and the metric by

$$126 \quad d(T_1, T_2) = \sum_{0 \leq i, j \leq n} \mu(c_1^{(i,j)}, c_2^{(i,j)})$$

127 This metric is not only defined for trees of the same size, but also for trees of  
128 different sizes where it’s natural to assign a coefficient of 0 to each term that is absent in a  
129 polynomial.

*Parameter estimation and model selection*

130

131 To estimate parameters for trees, we use the polynomial metric or the  
132 subtree-Manhattan metric together with the weighted average of the neighboring observed  
133 data with the nearest neighbor kernel smoother. Specifically, we generate a set of observed  
134 trees  $\mathcal{T}$  using a random tree generator with the different vectors of parameters  $\rho$ . For any  
135 tree  $T$  in  $\mathcal{T}$ , let  $\rho(T)$  be the vector of parameters used to generate  $T$ . We estimate the  
136 parameters of a tree  $T_0$  by the weighted average as follows:

$$137 \quad \hat{\rho}(T_0) = \frac{\sum_{T \in \mathcal{T}} K(T_0, T) \rho(T)}{\sum_{T \in \mathcal{T}} \rho(T)}$$

138 where  $K(T_0, T)$  is the  $k$ -nearest-neighbor kernel function, that is,  $K(T_0, T) = 1/k$  if  $T$  is a  
139  $k$  nearest neighbor of  $T_0$  under the polynomial metric and  $K(T_0, T) = 0$  otherwise. We  
140 choose different  $k$  for different sets of observed trees. For a set of observed trees  $\mathcal{T}$ , we  
141 generate another set  $\mathcal{S}$  of 1,000 random trees. For each  $k$  from 1 to 20, we estimate  
142 parameters of trees in  $\mathcal{S}$  using the set of observed trees  $\mathcal{T}$ , and we have the average  
143 estimation error for each  $k$ . We choose the  $k$  that has minimum average estimation error  
144 for the set of observed trees  $\mathcal{T}$ .

145 We use naive Bayes classifiers (Rish, 2001) together with the polynomial to perform  
146 model selection. Naive Bayes classifiers assume independence of the predictor variable. We  
147 label each tree according to the underlying model (beta splitting, the explosive radiation  
148 and trait evolution), and use the trees' polynomial coefficients as features.

### 149 *Simulations*

150 *Beta splitting trees* The beta splitting random trees used in this paper are  
151 generated by the beta-splitting model introduced in (Aldous, 1996). At each branching  
152 event, the probability of one child clade containing  $i$  tips and the other child clade  
153 containing  $n - i$  tips is given by the following formula.

$$154 \quad p(i|n) = \frac{1}{a_n(\beta)} \frac{\Gamma(\beta + i + 1)\Gamma(\beta + n - i + 1)}{\Gamma(i + 1)\Gamma(n - i + 1)}$$

155 The  $\Gamma(z)$  in the formula is the Gamma function and  $a_n(\beta)$  is a normalizing constant.

156 Our sets of  $n$ -tip modeled beta splitting trees consist of trees generated with  $\beta = 0$ ,  
157  $\beta = -1$ , and  $\beta = -1.5$ , and there are 100 trees for each parameter. These choices of  $\beta$   
158 correspond to the Yule model, the Aldous branching model and the proportional to  
159 distinguishable arrangements (PDA) model (Blum, 2006). We also use sets of beta  
160 splitting trees consisting of 1,000 such trees, with  $n$  tips and parameters  $\beta$  that are  
161 uniformly randomly chosen from the interval  $[-1.5, 8.5]$ .



162 *Explosive radiation trees* The explosive radiation trees were simulated with a  
163 modification of the birth-death model proposed by Steel (2001). Steel’s model builds on  
164 the traditional constant birth-death model by setting lineage-specific speciation rates.  
165 More precisely, the rate of speciation events on a given lineage is a function of  $t$ , the time  
166 to the last speciation event on that lineage. This time  $t$  is reset to 0 at every speciation,  
167 and the birth ( $\lambda_i$ ) and death ( $\mu_i$ ) rates of a given lineage  $i$  are then defined as follows:

$$168 \lambda_i(t) = \begin{cases} \lambda_B & \text{if } t < \tau \\ \lambda_A & \text{otherwise} \end{cases}$$
$$\mu_i(t) = \mu,$$

169 where  $\lambda_A$ ,  $\lambda_B$ ,  $\mu$  and  $\tau$  are parameters of the model.

170 All rates are defined as the number of events per tip per time unit. The choice of  
171 the time unit is not relevant to our experiments, as the polynomial does not make use of  
172 information on branch lengths.

173 A data set of  $n$ -tip explosive radiation trees contains 1,000 random trees generated  
174 with the birth rate  $\lambda_B$  fixed at 1.0 (per arbitrary time unit), the time shifting the birth  
175 rates  $\tau$  fixed at 0.5 time unit, and both the birth rate  $\lambda_A$  and the death rate  $\mu$  uniformly  
176 randomly chosen from the interval  $[0, 1]$ .

177 *Trait evolution trees* This data set was simulated following the birth-death model  
178 proposed by Heard (1996). In this model, each lineage has an associated trait value ( $x$ )  
179 which is “inherited” at speciation events with some stochastic change. The model for trait  
180 evolution implemented here is a linear-Brownian variation, where additive changes are  
181 made to the trait value at each speciation event:  $x_{\text{new}} = \max\{x_{\text{old}} + \epsilon, 0.01\}$ . The stochastic  
182 change  $\epsilon$  is drawn from a normal distribution with expectation zero and standard deviation  
183  $\sigma_x$ . Both  $\sigma_x$  and  $x_0$  (the trait value at the root) are parameters of the model.

184 The birth ( $\lambda_i$ ) and death ( $\mu_i$ ) rates are defined as  $\lambda_i = x$  and  $\mu_i = \mu$ , respectively.  
185 Similar to the explosive radiation model, the death rate  $\mu$  is constant in time and across  
186 lineages. Notice that there are numerous ways to produce trees with a given number of

187 species from an evolutionary model (Hartmann, 2010). For all evolutionary models used in  
188 our analysis, trees are simulated forward in time until  $n$  tips are first reached. Our data  
189 sets of  $n$ -tip trait evolution trees contain 1,000 random trees generated with the initial  
190 birth rate fixed at 1.0 (per arbitrary time unit), and the birth rate variation at a speciation  
191 event and the death rate uniformly randomly chosen from the interval  $[0, 1]$ .

192 We do not down-sample the simulated trees despite the fact that the data we use  
193 (see below) are only a small minority of the true numbers of tips in the relevant settings.  
194 This would be infeasible at genuine scales given the comparatively high true population  
195 sizes of circulating pathogens. For example, only a very small minority of circulating  
196 influenza infections lead to a sequence deposition in the database, with many others going  
197 undetected and/or unsequenced. Those that are sequenced may not be unique exemplars  
198 of their sequences in the population, as transmission may occur without detectable  
199 variation. As a consequence, in comparing simulation models to data, we interpret  
200 simulated branching events as diversification events that are likely to be ancestral to  
201 sampled tips and therefore observed, and "death" events as, effectively, sampling events  
202 that stop onward transmission of the particular lineage.

### 203 *Data*

204 *HIV and influenza virus trees* The HIV trees were described and analyzed  
205 previously (Chindelevitch, 2019). Briefly, HIV-1 sequence data from three studies were  
206 used. The Wolf et al. study (Wolf, 2017) provided data from a concentrated epidemic of  
207 HIV-1 subtype B, occurring primarily in men who have sex with men (MSM) in Seattle,  
208 USA. The Novitsky et al. study (Novitsky, 2013) describes data from a generalized  
209 epidemic of HIV-1 subtype C in Mochudi, Botswana, a village in which the HIV-1  
210 prevalence in the adult population at the time was estimated to be approximately 20%.  
211 Hunt et al. (Hunt, 2013) describes data from a national survey of the generalized epidemic  
212 of HIV-1 subtype C in South Africa. These datasets reflect a diverse set of spatial scales

213 and epidemiological contexts. Phylogenetic reconstruction was described in (Chindelevitch,  
214 2019); briefly, trees were reconstructed using RAxML (Stamatakis, 2014), which is a  
215 maximum likelihood method, under a general time-reversible (GTR) model of nucleotide  
216 substitution. We use a GTRCAT model for rate variation among sites. Each tree was  
217 based on a random sample of 100 sequences. We use a subtype D sequence as an outgroup  
218 to root HIV-1 subtype B phylogenies.

219 Our influenza virus trees were previously described in (Colijn, 2018). We aligned  
220 HA protein sequences from NCBI, focusing on human influenza A (H3N2). Data were  
221 downloaded from NCBI on 22 Jan. 2016. We included full-length HA sequences with  
222 collection date. The USA dataset ( $n = 2168$ ) includes sequences from the USA with  
223 collection dates between Mar. 2010 and Sep. 2015. The tropical dataset ( $n = 1388$ )  
224 includes sequences with a location listed as tropical, with collection dates within Jan. 2000  
225 and Oct. 2015. Accession numbers are included in the Supporting Information of Colijn  
226 (2018). Fasta files were aligned with mafft, and for both the tropical and USA datasets,  
227 500 taxa were selected uniformly at random 200 times. We inferred 200 corresponding  
228 phylogenetic trees with FastTree (Price, 2010). Where necessary we re-aligned the 500  
229 selected sequences before performing tree inference. This process resulted in 200 “tropical”  
230 influenza virus trees and 200 “USA” influenza virus trees, each with 500 tips,  
231 reconstructed from the HA region of human H3N2 samples. Note that this approach is  
232 distinct from the perhaps more familiar phylogenetic methods where bootstrapping or  
233 Bayesian reconstructions results in many trees on *one* set of tips. These are likely to share  
234 features and structures because they describe the ancestry of the same set of taxa. Here,  
235 each tree has a different set of tips (though there is some overlap).

236 *WHO influenza virus clades* We used several influenza virus clades, described in  
237 (Hayati, 2020). In that work we downloaded all human H3N2 full-length HA sequences  
238 with dates between 1980 and May 2018 and created a large, timed phylogeny of H3N2  
239 using RAxML and Least Squares Dating (Stamatakis, 2014; To, 2016). This “full” tree has

Data	Data source	Virus	Size (tips)
Wolf (100 trees)	Wolf (2017)	HIV-1 subtype B	500
Novitsky (100 trees)	Novitsky (2013)	HIV-1 subtype C	500
Hunt (100 trees)	Hunt (2013)	HIV-1 subtype C	500
USA (200 trees)	NCBI	Human influenza A	500
Tropical (200 trees)	NCBI	Human influenza A	500
A1B/135N	NCBI	Human influenza A	60
A1B/135K	NCBI	Human influenza A	63
3c3.B	NCBI	Human influenza A	117
A3	NCBI	Human influenza A	227

Table 1. Summary of virus phylogenies.

240 over 12,000 tips. We used the Nextflu (Neher, 2015) *augur* pipeline  
241 (<https://bedford.io/projects/nextflu/augur/>) to assign a WHO clade designation to  
242 the sequences. The WHO defines named clades using specific mutations in the HA1 and  
243 HA2 subunits of the HA protein. The full list of mutations is available at: [https://github.com/nextstrain/seasonal-flu/blob/master/config/clades\\_h3n2\\_ha.tsv](https://github.com/nextstrain/seasonal-flu/blob/master/config/clades_h3n2_ha.tsv).  
244 We assign a sequence to a clade if it contains all the mutations defining that clade. We  
245 then extracted the subtrees of the “full” tree corresponding to specific WHO clades  
246 A1B/135N (60 tips), A1B/135K (63 tips), 3c3.B (117 tips) and A3 (227 tips). These are  
247 recent and appropriately-sized trees which we use here to demonstrate parameter  
248 estimation for simple models, and model selection among our four random tree models.  
249

### 250 *Implementation*

251 We developed an R package named *treenomial*, which is available at CRAN. We  
252 also prepared a demonstration named *treeverse*, which displays a 3-dimensional projection  
253 of the polynomial metric space of all binary tree shapes up to 16 tips with interactive  
254 options available at <https://magpiegroup.shinyapps.io/treeverse/>.

RESULTS

*Tree Representations and Metrics*

We compare the polynomial to other tree representation methods in terms of computation time and resolution. These tree-representing methods include the Colless index (Agapow, 2002), gamma statistics (Pybus, 2000), the Sackin index (Sackin, 1972), the subtree size distribution and more recently introduced Laplacian spectrum (Lewitus, 2016). The resolution of a tree-representing method (for  $n$ -tip trees) is defined to be the ratio of the number of unique representations to the total number of non-isomorphic tree shapes with  $n$  tips. We compute these representations for all tree shapes with 15 tips (where there are 4850 non-isomorphic tree shapes). Figure 1 A displays the computation time and the resolution of these methods, where the data point “combined” is the vector comprising the Colless index, gamma statistics and the Sackin index. The results show that Laplacian spectrum, the polynomial, and the subtree size distribution (with more than one parameter) have higher resolution than scalar summary statistics while the scalar Colless index, gamma statistics and the Sackin index have lower resolution. As there are vastly numerous non-isomorphic tree shapes with hundreds of tips, it is not feasible to compute the resolution for larger trees, but we know that the resolution of the subtree size distribution decreases as the number of tips increases, and the Laplacian spectrum is not guaranteed to have 100% resolution for all trees, that is, there are non-isomorphic trees with the same spectrum density (Lewitus, 2016). The polynomial, on the other hand, is guaranteed to distinguish all trees (Liu, 2021). In Figure 1 B, we show how computation time of the subtree size distribution, the Laplacian spectrum and the polynomial for a single tree changes as the size of trees increase. Among the high-resolution tree-representing methods we compared, the polynomial has low computation time and keeps the resolution at 100% for trees of any size.

Tree representations can induce tree metrics, which are important tools in comparing phylogenetic trees. We compare the polynomial metric with the metric induced

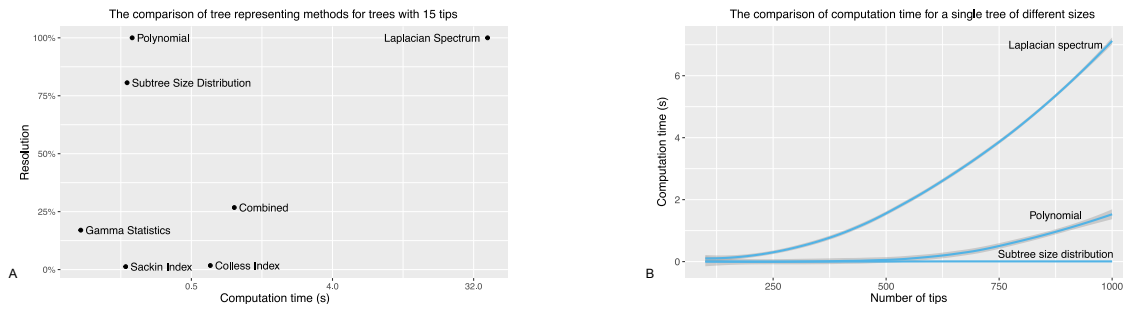


Figure 1. A: the comparison of tree representing methods, where the combined is the a combined vector of the Colless index, gamma statistics and the Sackin index. B: the comparison of the computation time for trees of different sizes. These are based on the computation time for random trees with 100 to 1,000 tips with increment of 50 tips; each data point denotes the average computation time over 1,000 trees.

282 by high-resolution tree-representing methods, that is, Lewitus-Morlon metric and the  
283 subtree-Manhattan metric. The polynomial metric is a genuine metric on trees, in the  
284 sense that it only gives a distance of zero if two trees have identical shapes, it is  
285 symmetric, and obeys the triangle inequality (see the supplement for proof; in contrast, the  
286 subtree-Manhattan metric and Lewitus-Morlon metric are not metrics in the mathematical  
287 distance sense (Lewitus, 2016)). The polynomial metric also has the advantage that the  
288 distance between a pair of trees is bounded above by the number of non-zero entries in the  
289 coefficient matrix of the larger tree. More precisely, let the larger tree be of  $n$  tips; the  
290 polynomial distance between the trees has an upper bound of  $n \lfloor n/2 \rfloor - \lfloor n/2 \rfloor^2$ . The  
291 distribution of the pairwise distances between trees of the same size resembles a normal  
292 distribution, which gives a relative reference for how large the distance between a pair of  
293 trees is compared to what one might expect. See Supplementary Figure 1 for the  
294 distribution. Figure 2 A-C displays visualizations of the three distances between trees in a  
295 data set of 100-tip modeled beta splitting trees. We apply the  $k$ -medoids clustering  
296 algorithm PAM described in (Kauffman, 1990) to, respectively, the Lewitus-Morlon  
297 distance matrix, the subtree-Manhattan distance matrix and the polynomial distance  
298 matrix of a set of the 100-tip modeled beta splitting trees. We repeat this experiment for  
299 100 times; Figure 2 D shows the misclassification rates. The polynomial metric has smaller  
300 misclassification rates than the other two metrics, which indicates that the polynomial has

301 the potential to perform better in tasks involving clustering phylogenetic trees.

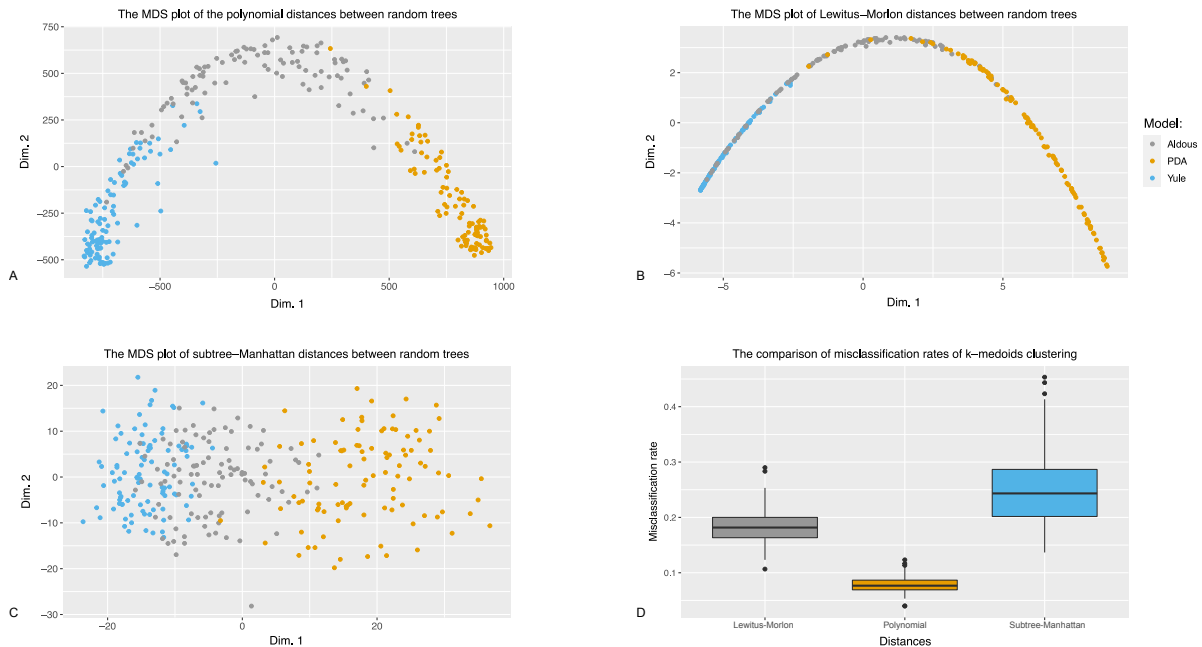


Figure 2. A–C: the multidimensional scaling plots of the three distances between trees in a set of 100-tip random trees, where each dot represent a random tree. D: the comparison of the misclassification rates of  $k$ -medoids clustering.

302

### *Parameter Estimation and Model Selection*

303

*Parameter estimation* We show that the polynomial can be used to create

304

likelihood free methods for parameter estimation. Here, we display the results of parameter

305

estimation using the polynomial metric together with a simple weighted average method

306

described in the method section. We generate a set of 250-tip beta splitting trees and use

307

the set of random trees as observed data in the parameter estimation method; we then

308

estimate the parameter  $\beta$  for 100 beta splitting trees with 250 tips. Figure 3 A shows the

309

result of the estimation, and Figure 3 B shows the result of the estimation for 500-tip beta

310

splitting trees. See Supplementary Table 1 for the summary of the estimation. In general

311

the estimation works better for larger trees, and is better when the parameter  $\beta$  is in the

312

interval  $[-1.5, 2]$ . We note that where a likelihood model is available, maximizing the

313 likelihood may well be better than likelihood-free inference based on tree descriptions, but  
314 these results indicate that the polynomial contains relevant information that  
315 high-performance likelihood-free inference methods could utilize.

316 We also generate a set of 750-tip explosive radiation trees and use the set of random  
317 trees as observed data in the parameter estimation method to estimate the birth rate  $\lambda_B$   
318 and death rate  $\mu$  for 100 explosive radiation trees with 750 tips. Figure 3 C-D shows the  
319 results of the estimation. The results are not as good as the results for beta splitting trees,  
320 especially the results for the birth rate  $\lambda_B$ . Supplementary Tables 1-3 give details of the  
321 relationship between estimated and true values. We also use the subtree-Manhattan metric  
322 and the same weighted average method to perform parameter estimation for the same data  
323 sets. See Supplementary Figure 4; we find that the polynomial metric performs better than  
324 the subtree-Manhattan metric with the weighted average method in estimating parameters  
325 for both beta splitting trees and the explosive radiation trees.

326 *Model selection* The beta splitting model and the explosive radiation model are  
327 different random tree generators. The beta splitting model uses the Markov branching  
328 process while the explosive radiation model uses a birth-death process. Both processes are  
329 commonly used in random tree generators, for example, the trait evolution model is  
330 another tree generator based on the birth-death process. Figure 2 shows that the  
331 polynomial has the potential to distinguish different tree generating models. In this  
332 section, we use the polynomial together with naive Bayes classifiers to estimate which  
333 model is used to generate a tree.

334 We generate a set of 500-tip beta splitting trees, a set of 500-tip explosive radiation  
335 trees, and a set of 500-tip trait evolution trees. We use these sets of random trees as  
336 observed data together with the naive Bayes classifiers to classify random trees generated  
337 by these three models. Figure 3 E shows the results of the experiment where we only use  
338 the set of beta splitting trees and the set of explosive radiation trees to train the naive  
339 Bayes classifier, and use the classifier to classify a set of 1,000 beta splitting trees and 1,000



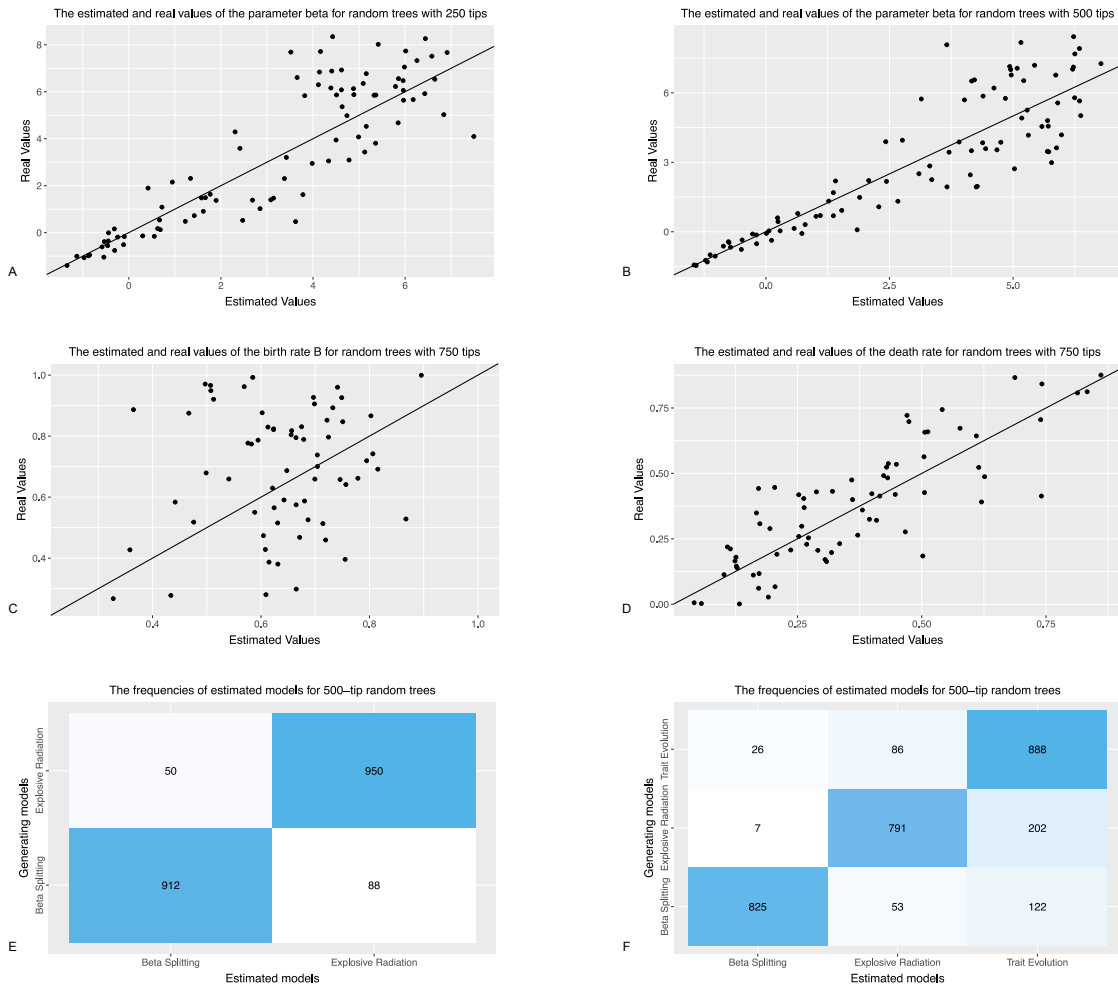


Figure 3. A-B: the comparisons between the real parameter and the estimated parameter of the beta splitting random trees with 250 tips and 500 tips using polynomials. C-D: the comparisons between the real parameters and the estimated parameters of the explosive radiation random trees with 750 tips using polynomials. E-F: the results of using naive Bayes classifiers to select the model generating random trees with 500 tips using polynomials.

340 explosive radiation trees. Figure 3 F shows the results of the experiment where we use all  
 341 three sets of random trees to train the naive Bayes classifier and use the classifier to classify  
 342 a set of 1,000 beta splitting trees, 1,000 explosive radiation trees, and 1,000 trait evolution  
 343 trees. The accuracy of the first experiment is 93.1% and of the second experiment is 83.5%,  
 344 where the main misclassification (58.1% of the misclassified cases) is between the explosive  
 345 radiation model and the trait evolution model, the two models based on the birth-death  
 346 processes. Supplementary Figure 3 shows the results for 250-tip and 750-tip trees, and that

347 this model selection method is more robust with larger trees. The results show that the  
348 polynomial together with naive Bayes classifiers can be a good tool in finding a tree  
349 generator that fits a given tree, as not only are trees from different random processes  
350 distinguished well, but the two different birth-death processes are also well distinguished.

351 We also use the subtree size distribution and the standard naive Bayes classifiers to  
352 perform model selection for the same data sets. See Supplementary Figure 4. Compared to  
353 the polynomial, the accuracy of the first experiment using the subtree size distribution is  
354 82.7% and of the second experiment is 71.6%, where more explosive radiation trees are  
355 classified as beta splitting trees. To further understand the differences between the  
356 polynomial and the subtree size distribution in the naive Bayes classifiers, we display the  
357 most informative features in the classifiers in Supplementary Figure 5. It shows that for  
358 the subtree size distribution, the most informative features in model selection are the  
359 number of subtrees with approximately 400 tips, which could be considered as a  
360 description of tree imbalance for more imbalanced trees would have more subtrees with 400  
361 tips than the balanced ones. On the other hand, Supplementary Figure 5 B shows that  
362 other than the clade size distribution (the dark thin strip at the bottom), the most  
363 informative features for the polynomial also include the coefficients in the black area at the  
364 top, which are interpreted as the numbers of ways to choose as many clades (with more  
365 than one tips) as possible so that the clades contain all or most of the tips in a tree.  
366 Compared to the subtree size distribution, this extra information gives the performance of  
367 the model selection method a boost.

### 368 *Applications to Data*

369 Human influenza virus A H3N2. Influenza virus A is highly seasonal outside the  
370 tropics and most cases occur in the winter (Russell, 2008), whereas there is relatively little  
371 seasonal variation in the tropics. This demonstrative data set provides trees for the same  
372 virus circulating with different epidemiological dynamics (seasonal forcing in temperate

373 regions, vs lack of seasonality in the tropics). The second data set consists of three samples  
 374 of trees inferred from HIV-1 sequences in different settings: subtype B among men who  
 375 have sex with men (MSM) in Seattle (Wolf, 2017), HIV 1C circulating at the village scale  
 376 in Botswana (Novitsky, 2013) and a national-level dataset from South Africa (Hunt, 2013).  
 377 As with influenza virus, it is to be hoped that these different epidemiological patterns are  
 378 revealed in the shapes of reconstructed phylogenetic trees (Chindelevitch, 2019; Colijn,  
 379 2018).

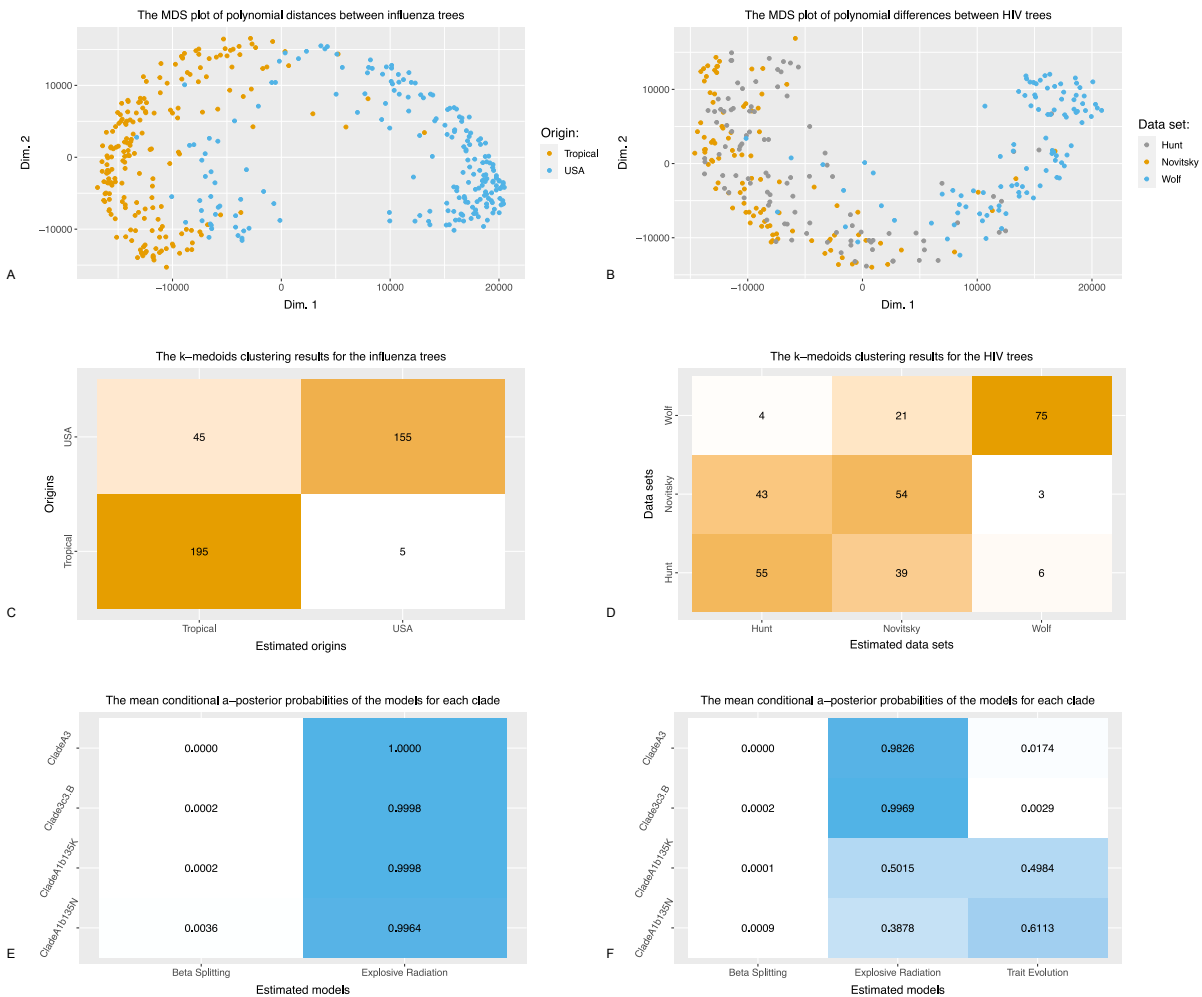


Figure 4. A: the MDS plots of the polynomial distances between the influenza trees. B: the MDS plots of the polynomial distances between the HIV trees. C: the results of  $k$ -medoids clustering for the influenza trees using the polynomial metric. D: the results of  $k$ -medoids clustering for the HIV trees using the polynomial metric. E-F: the mean conditional a posteriori probabilities (over the 1,000 naive Bayes classifiers) of the model estimation for the influenza clades.

380 We visualize the polynomial distances between trees in these two sets by classical  
381 MDS in Figure 4 A-B. We also use the  $k$ -medoids clustering on the data and we have the  
382 results displayed in Figure 4 C-D. The influenza trees are very well separated into desired  
383 groups under the  $k$ -medoids clustering. This indicates that classifying the epidemiological  
384 process behind a tree using the polynomial metric would likely be possible. In the  
385 supplement, we also compute the binary differences (Choi, 2010) of the polynomials for  
386 these trees, which improves the results of the  $k$ -medoids clustering. See Supplementary  
387 Figure 6. For these particular challenges, however, typically a researcher would know  
388 whether their data were from the tropics or not, or what the broad epidemiological setting  
389 (village, national, Western population MSM) was at the time of collection. We therefore  
390 focus on more specific estimation questions (parameter estimation and model choice).

391 As an example of applying the parameter estimation and model selection methods  
392 to data, we first select the models that best fit the four WHO influenza clades, A1B/135N  
393 (60 tips), A1B/135K (63 tips), 3c3.B (117 tips) and A3 (227 tips), then estimate the  
394 parameters for the model that best fits the clade. To select the model that best fits a clade,  
395 we generate a set of beta splitting trees, a set of explosive radiation trees and a set of  
396 trait evolution trees of the same size as the clade. We use these sets of trees and naive  
397 Bayes classifiers to estimate the a-posterior probabilities of the clade being generated by  
398 the models. Figure 4 E shows that if we select only from the beta splitting model and the  
399 explosive radiation model, then all four clades are deemed more likely to be generated by  
400 the explosive radiation model, a tree generator based on the birth-death model. Figure 4 F  
401 shows that if we include the trait evolution model, the small clades A1B/135N (60 tips)  
402 and A1B/135K (63 tips) are predicted to be generated by either the explosive radiation  
403 model or the trait evolution model. The classifiers predict that for larger clades, the most  
404 likely model is the explosive radiation model. Both models seem reasonable for influenza,  
405 as a new variant that has polymorphisms allowing it to evade immunity that has built up  
406 in the population due to exposure to previous influenza viruses could have an early rapid

407 rise (explosive radiation). However, influenza viruses have numerous traits (including  
408 interactions with host immunity) that could influence the branching rates in influenza  
409 virus phylogenies.

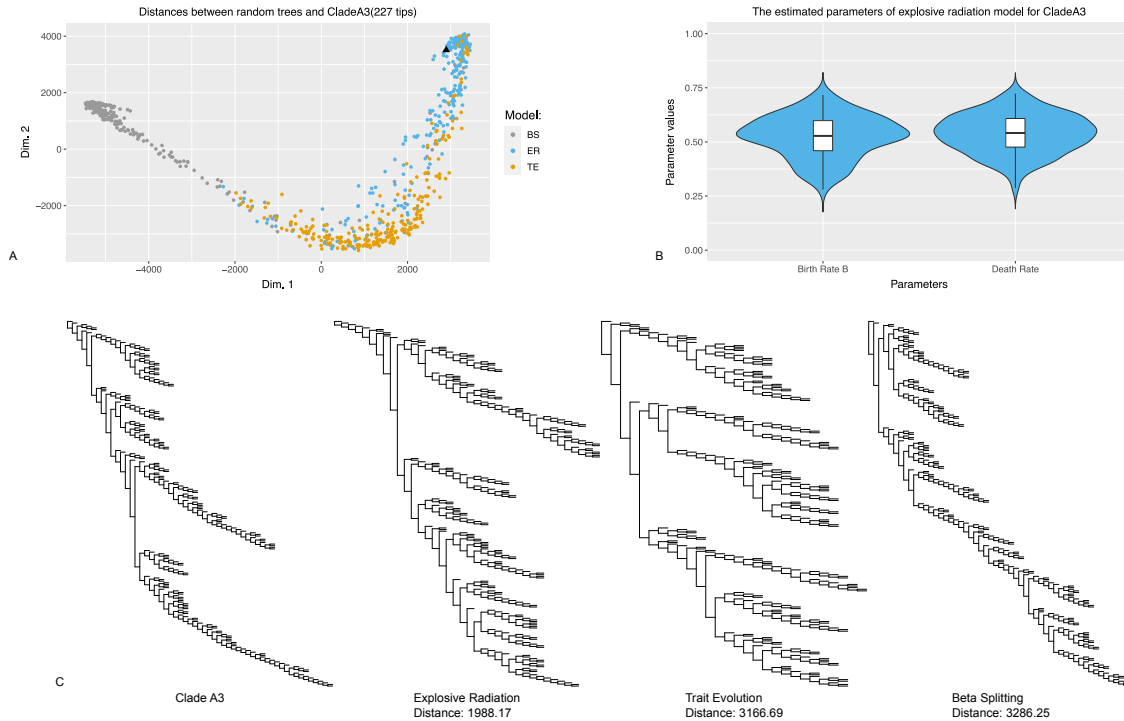


Figure 5. A: The MDS plots of the polynomial distances between the random trees generated by the three different models and influenza virus clade A3. B: the distribution of the estimated parameters of clade A3 over 100 replicates. C: the plots of clade A3 and the nearest random trees generated by the three different models.

410 As an example, we examine influenza virus clade A3 (227 tips) in detail and  
411 estimate its parameters. First, we generate a set of beta splitting trees, a set of explosive  
412 radiation trees, and a set of trait evolution trees, all with 227 tips. For each set, we choose  
413 250 random trees to visualize. In total we thus have 751 trees including clade A3. Figure  
414 5A shows a visualization of the polynomial distances among these trees. Figure 5B shows  
415 the results of estimating the parameters (repeated 100 times with different sets of random  
416 trees) of the explosive radiation model for clade A3. The 95% confidence interval of the  
417 birth rate  $\lambda_B$  is (0.50, 0.54) and the 95% confidence interval of the death rate  $\mu$  is  
418 (0.52, 0.56). The 95% confidence interval of  $R_0$  of the clade is (0.906, 1.019).

419 Lastly, we plot, in Figure 5 C, clade A3 and the nearest random trees to the clade  
420 from each model in the sets of 250 random trees displayed in Figure 5 A. The polynomial  
421 distances between a pair of 227-tip trees has the upper bound of 12,882. Assuming the  
422 distribution of the pairwise distances will be normal as displayed in Supplementary  
423 Figure 1, all of the polynomial distances between clade A3 and the trees in Figure 5 C are  
424 below average pairwise distances of 227-tip trees. We also perform the same analysis on  
425 clade 3c3.B which has 117 tips (Supplementary Figure 7). Throughout our comparisons  
426 between simulated and real trees, we note that we have not simulated realistic total  
427 populations of either HIV or influenza in the relevant settings and then down-sampled to  
428 match the sizes of observed trees as this would be infeasible. This affects the interpretation  
429 of our estimates.

## 430 DISCUSSION

431 We have introduced a new way to describe and analyze phylogenetic trees using a  
432 polynomial that uniquely characterizes trees. We compare the polynomial to other indices  
433 and methods describing tree shapes. The polynomial is easy to compute and it has the  
434 advantage of describing trees in full resolution, that is, the descriptions are different if and  
435 only if the two tree shapes are not isomorphic. Moreover, the polynomials have the  
436 potential to be extended to record information about the branch lengths.

437 We also introduced some basic methods for tree analysis using the polynomial. The  
438 methods discussed in this paper include a tree metric, a parameter estimation method  
439 based on the metric, and a naive Bayes classifier directly trained by the coefficients of the  
440 polynomials. We chose these simple and tractable methods to show that the polynomial  
441 can be utilized in likelihood free methods for various tasks in analyzing phylogenetic trees.  
442 These polynomial based methods can distinguish trees from different models and different  
443 data sets, help estimate parameters, and aid in model selection. We have also applied these  
444 polynomial based methods to estimate parameters and select the best-fit model for the

445 chosen WHO influenza virus clades. The results show that the tree shapes of the influenza  
446 clades are most similar to random trees generated by either the explosive radiation model  
447 or the trait evolution model, both of which are based on the birth-death process compared  
448 to the beta splitting model which is based on the Markov branching process. We also  
449 computed the nearest (in the polynomial distance) trees from each model to a  
450 WHO-defined influenza clade. This information, together with the distribution of the  
451 pairwise polynomial distances between trees being normal, can be used to assess how well  
452 a simulated tree resembles a tree reconstructed from data.

453 The simple methods used in this paper for parameter estimation and model  
454 selection can be improved in terms of computation efficiency among other aspects. And  
455 indeed, in estimation problems, it may be best to collect a wide range of tree descriptors  
456 (including polynomial coefficients, scalar summaries such as the Sackin and Colless  
457 imbalance measures and other high-resolution characterizations of the tree) (Saulnier,  
458 2017), and let feature selection sort out which are best for a particular problem. Different  
459 models and data will yield trees with different features, and in some of these, simple scalar  
460 summary statistics may perform well. Our results show that in our simulation examples  
461 the polynomial coefficients are informative and would likely add to such an analysis,  
462 probably with the most benefit where scalar imbalance measures do not contain sufficient  
463 information about trees to perform the desired estimation task. Characterizing trees in the  
464 polynomial's high-resolution metric way also allows selection of the closest tree to a tree  
465 from data, and visualizations of the space of trees derived from a model or datasets of  
466 interest. The polynomial can be used to obtain novel features or pseudo-metrics for  
467 clustering and estimation; as an example, the binary differences (Choi, 2010) can be used  
468 to improve clustering for the influenza and HIV trees (Supplementary Figure 6).

469 Our polynomial is not the only one that uniquely represents rooted binary trees.  
470 Other polynomials, such as the ones introduced in (Andr en, 2009), (Chaudhary, 1991),  
471 (Negami, 1996) and (Botti, 1993), (Matsen, 2012) are also good candidates for tree

472 analysis. Thus it is worth investigating more about how these different polynomials can be  
473 used to analyze phylogenetic trees and how different results can they yield. The  
474 computation of most of these polynomials requires going through all subtrees or all  
475 permutations of a given size, which can be computationally heavy, while the polynomial  
476 used in this paper has a recursive formula which makes the computation more efficient.

477 To compare trees with different sizes is another challenge in tree comparison. In this  
478 paper, we have compared trees with the same number of tips and we have proposed a way  
479 to compare trees with different sizes. In our demonstration *treeverse*, trees with different  
480 sizes are compared and the distances between the trees are visualized by an interactive 3-D  
481 MDS plot. There are various ways to compare the coefficient vectors and compare trees  
482 with different sizes, but for trees whose sizes are drastically different, the sizes naturally  
483 remain a dominating factor in the resulting tree comparisons.

484 Because polynomial coefficients can be treated as vectors, and vectors give rise to  
485 metrics, there are alternative metrics that can be defined using tree polynomials (both  
486 those used here and others (Andr n, 2009; Chaudhary, 1991; Negami, 1996)). Once trees  
487 are encoded as vectors, a range of regression, inference and dimension reduction and other  
488 machine learning tools can, as a result, be applied to trees. In addition, other tree shape  
489 statistics or further information about the trees (including measures of branch length) can  
490 easily be appended to the vectors to integrate distinct sources of data. This provides a  
491 scheme to study phylogenetic trees comprehensively.

492 There remains considerable scope to improve the clustering and classification tools  
493 used here, which we used to demonstrate that parameter estimation and model choice can  
494 be done. One challenge in this work is that there are too many polynomial coefficients;  
495 however, feature selection, hyperparameter optimization and dimension reduction tools  
496 could be used to reduce the number of features in a systematic way. Furthermore, we used  
497 one- or two-dimensional estimation tasks as demonstrations. Realistic models of evolution  
498 are likely to contain multiple parameters (for example, time-dependent speciation and



499 extinction rates; intra- and inter-group competition parameters, relative fitness), so more  
500 advanced and modern statistical inference tools could be considered. The simpler  
501 estimation we have provided is a proof of principle for using polynomial coefficients in  
502 estimation tasks.

#### 503 ACKNOWLEDGEMENTS

504 This work was supported by the grant of the Federal Government of Canada's  
505 Canada 150 Research Chair program to Dr. Caroline Colijn. We would like to thank Art  
506 Poon, who provided the HIV trees.

#### 507 SUPPLEMENTARY MATERIAL

508 *The polynomial metric* We prove that the polynomial metric is a genuine metric.  
509 It is easy to check that  $d(T_1, T_2) = 0$  if and only if  $T_1 \simeq T_2$ , and  $d(T_1, T_2) = d(T_2, T_1)$ . We  
510 show that the triangular inequality is true for the metric, that is,  
511  $d(T_1, T_3) \leq d(T_1, T_2) + d(T_2, T_3)$ . We only need to prove the following inequality holds for  
512  $a, b, c \geq 0$ .

$$513 \frac{|a - c|}{a + c} \leq \frac{|a - b|}{a + b} + \frac{|b - c|}{b + c}$$

514 Note that if  $a \geq c \geq b$  or  $c \geq a \geq b$ , we have

$$515 \frac{|a - c|}{a + c} \leq \frac{|a - b|}{a + c} + \frac{|b - c|}{a + c} \leq \frac{|a - b|}{a + b} + \frac{|b - c|}{b + c}$$

516 If  $a \geq b \geq c$ , we have

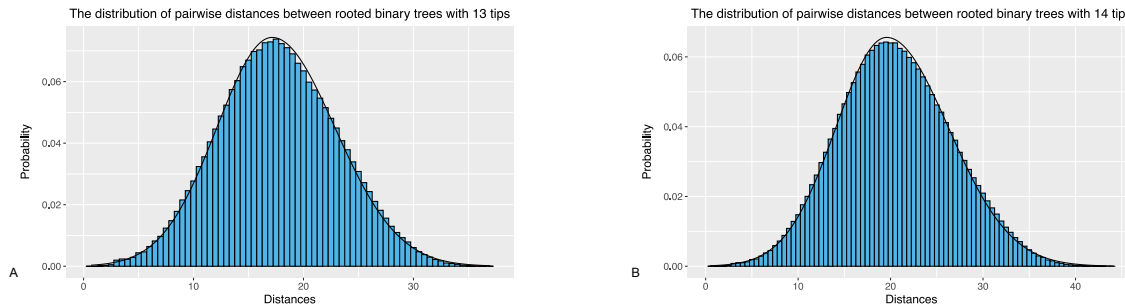
$$517 \frac{a - c}{a + c} \leq \frac{2b(a - c)}{(a + b)(b + c)} = \frac{a - b}{a + b} + \frac{b - c}{b + c}$$

518 This is equivalent to  $b^2 + ac \leq ab + bc$ , which is true because  $ac - bc \leq ab - b^2$ . Similarly,  
519 the equality also holds when  $c \geq b \geq a$ .

520 If  $b \geq a \geq c$ , we have

$$521 \frac{a - c}{a + c} \leq \frac{2(b^2 - ac)}{(a + b)(b + c)} = \frac{b - a}{a + b} + \frac{b - c}{b + c}$$

522 This is equivalent to  $ab(b - a) + 3c(b^2 - a^2) + c^2(b - a) \geq 0$ , which is true as  $b \geq a$ .  
523 Similarly, the equality also holds when  $b \geq c \geq a$ . Therefore the polynomial metric is a  
524 genuine metric.



Supplementary Figure 1. A-B: the distribution of all pairwise polynomial distances between all rooted binary trees with 13 and 14 tips, where the black solid curves are normal fits.

525 The distribution of polynomial distances between all pairs of trees with  $n$  tips  
526 resembles a normal distribution. Supplementary Figure 1 displays the distribution for trees  
527 with 13 and 14 tips, where the black solid curves are normal fits. For the distribution for  
528 13-tip trees, the estimated mean value is 17.70, the estimated standard deviation is 5.37,  
529 and Shapiro-Wilk normality test has  $W$  of 0.99 and p-value of  $6.21 \times 10^{-15}$ . For the  
530 distribution for 14-tip trees, the estimated mean value is 20.54, the estimated standard  
531 deviation is 6.10, and Shapiro-Wilk normality test has  $W$  of 0.99 and p-value of  
532  $4.43 \times 10^{-15}$ .

533 *Parameter estimation and model selection* We show the supplementary results of  
534 parameter estimation and model selection in complement to the figures displayed in the  
535 main result section. Supplementary Figure 2 shows the results of parameter estimation for  
536 750-tip beta splitting trees, 250-tip and 500-tip explosive radiation trees. Supplementary  
537 Table 1, 2 and 3 show the summaries of the estimation. Supplementary Figure 3 shows the  
538 results of model selection for 250-tip and 750-tip random trees generated by the three  
539 models. Supplementary Figure 4 shows the results parallel to the results displayed in  
540 Figure 3 with subtree size distributions instead of polynomials.

Models	$R^2$	Coeff.	Est.	Std. Err.	$p$ -value
BS-Poly 250 tips	0.77	Intercept	0.016	0.25	0.95
		Slope	1.06	0.06	<2e-16
BS-Poly 500 tips	0.78	Intercept	-0.038	0.22	0.86
		Slope	1.012	0.056	<2e-16
BS-Poly 750 tips	0.85	Intercept	0.026	0.21	0.23
		Slope	1.024	0.043	<2e-16
BS-SSD 250 tips	0.41	Intercept	-0.84	0.55	0.13
		Slope	1.009	0.13	4.31e-12
BS-SSD 500 tips	0.37	Intercept	-0.79	0.71	0.27
		Slope	1.076	0.16	9.52e-10

Supplementary Table 1. The summary of linear fit of the real parameter and the estimated parameter for beta splitting trees.

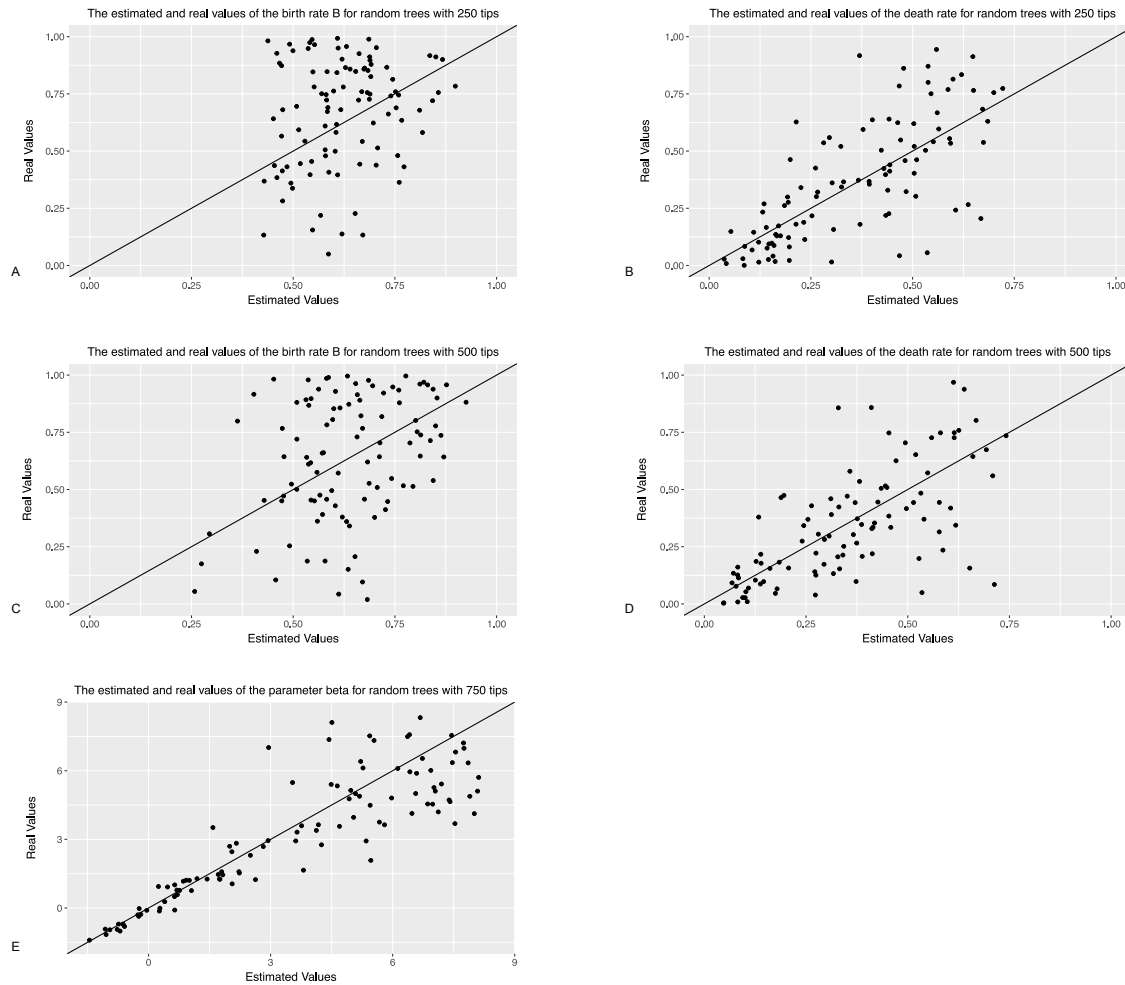
Models	$R^2$	Coeff.	Est.	Std. Err.	$p$ -value
ER-Poly 250 tips	0.035	Intercept	0.33	0.15	0.033
		Slope	0.46	0.25	0.063
ER-Poly 500 tips	0.043	Intercept	0.36	0.13	0.006
		Slope	0.43	0.21	0.037
ER-Poly 750 tips	0.09	Intercept	0.29	0.14	0.048
		Slope	0.58	0.22	0.011
ER-SSD 750 tips	0.07	Intercept	0.40	0.10	0.001
		Slope	0.44	0.16	0.008

Supplementary Table 2. The summary of linear fit of the real parameter and the estimated parameter  $\lambda_B$  for explosive radiation trees.

Models	$R^2$	Coeff.	Est.	Std. Err.	$p$ -value
ER-Poly 250 tips	0.24	Intercept	0.12	0.05	0.016
		Slope	0.67	0.12	1.87e-07
ER-Poly 500 tips	0.46	Intercept	0.067	0.041	0.11
		Slope	0.84	0.091	7.44e-15
ER-Poly 750 tips	0.69	Intercept	0.034	0.031	0.28
		Slope	0.94	0.075	<2e-16
ER-SSD 750 tips	0.59	Intercept	-0.008	0.035	0.815
		Slope	1.40	0.12	<2e-16

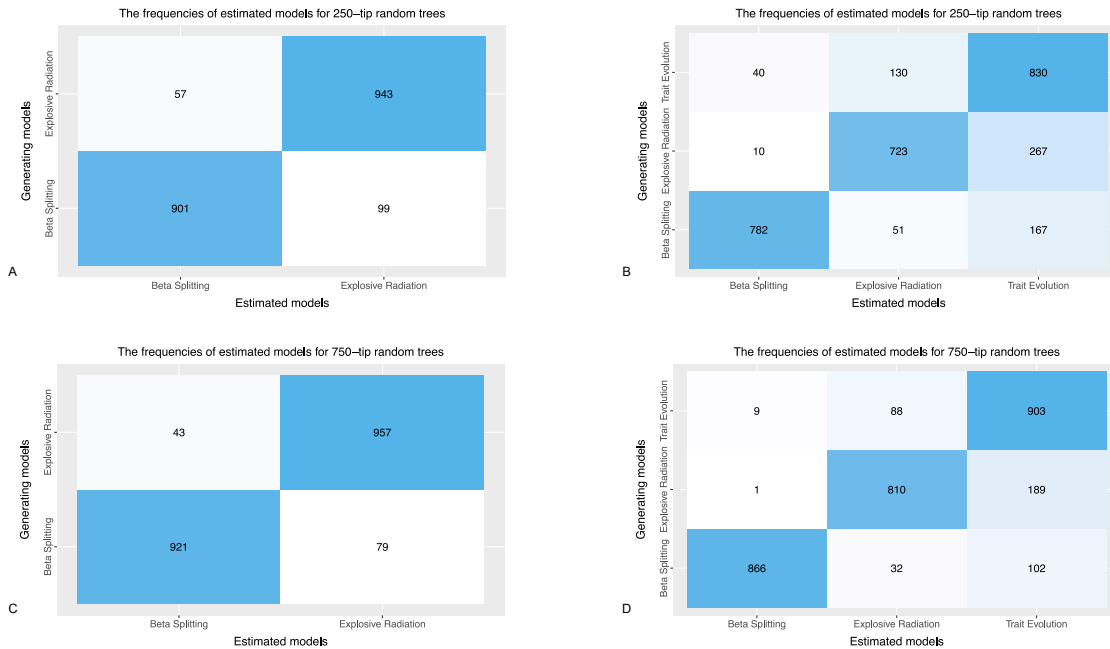
Supplementary Table 3. The summary of linear fit of the real parameter and the estimated parameter  $\mu$  for explosive radiation trees.

541 In Supplementary Figure 5, we show the importance of features in the naive Bayes  
 542 classifiers used for model selection with both subtree size distributions and polynomials.  
 543 As the naive Bayes classifiers assume independence of variables, the Shannon entropy  
 544 reflects the importance of the features, where a feature with smaller entropy means the  
 545 feature is more important in the classifier.



Supplementary Figure 2. A-B: the comparisons between the real parameters and the estimated parameters of the explosive radiation random trees with 250 tips using polynomials. C-D: the comparisons between the real parameters and the estimated parameters of the explosive radiation random trees with 500 tips using polynomials. E: the comparisons between the real parameter and the estimated parameter of the beta splitting random trees with 750 using polynomials.

546 *Polynomial binary differences* Binary differences, based on presence and absence  
547 of components, though in general not metrics, are one of the commonly used indices in, for  
548 example, taxonomic, ecologic, biogeographic comparison and classification (Choi, 2010).  
549 They provide effective insights about clusters though they are not metrics in general. We  
550 define the polynomial binary differences used in this paper by the number of terms that  
551 are present in the polynomial of one tree but are absent in the polynomial of the other.  
552 More precisely, the binary difference of two trees  $T_1$  and  $T_2$  are calculated by counting the



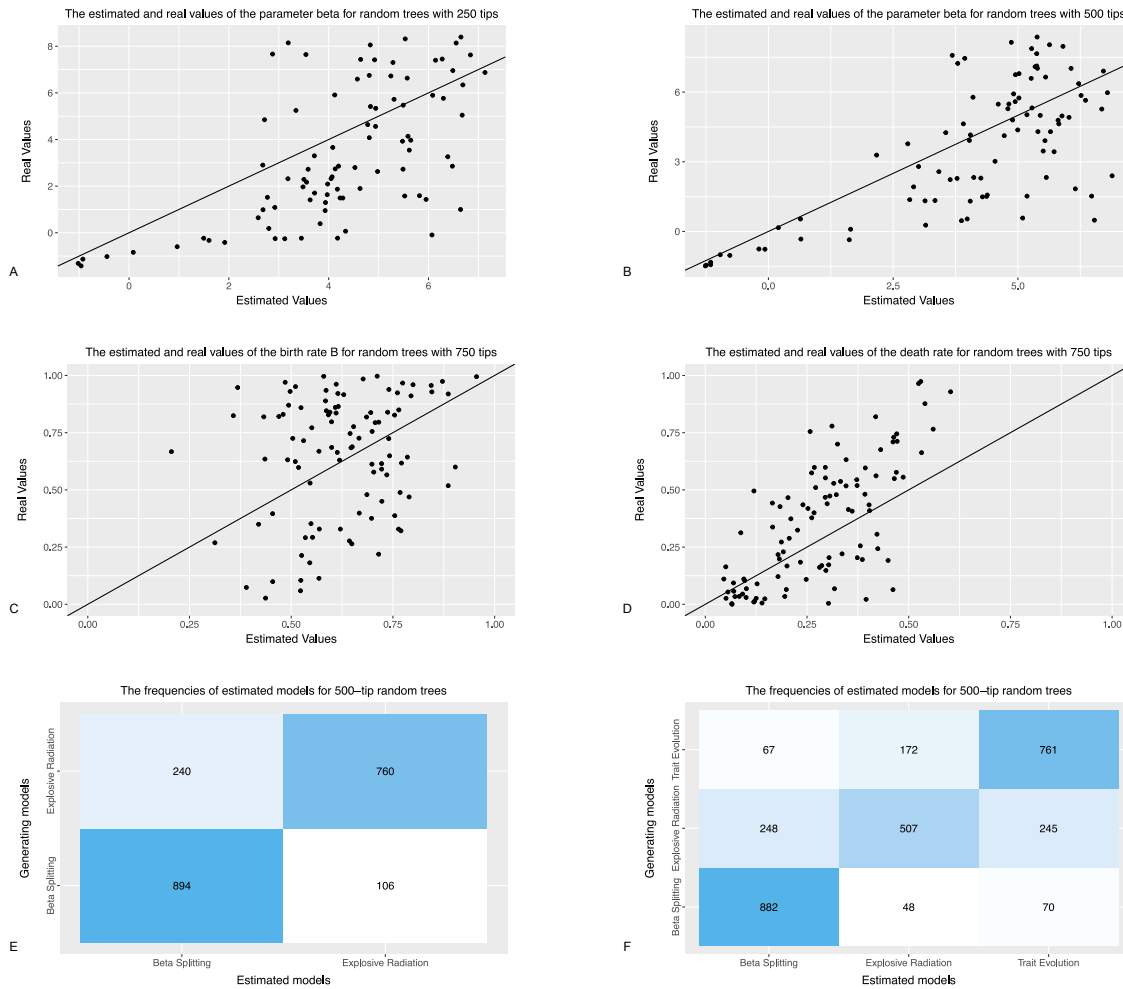
Supplementary Figure 3. A-B: the results of using naive Bayes classifiers to select the model generating random trees with 250 tips using polynomials. C-D: the results of using naive Bayes classifiers to select the model generating random trees with 750 tips using polynomials.

553 number of terms that are present in  $P(T_1, x, y)$  but are absent in  $P(T_2, x, y)$ , or the  
 554 number of terms that are present in  $P(T_2, x, y)$  but are absent in  $P(T_1, x, y)$ . This provides  
 555 another way to compare polynomials (trees). Supplementary Figure 6 shows the results of  
 556  $k$ -medoids clustering on the binary differences of the influenza trees and the HIV trees,  
 557 which are better than the polynomial metric in this task.

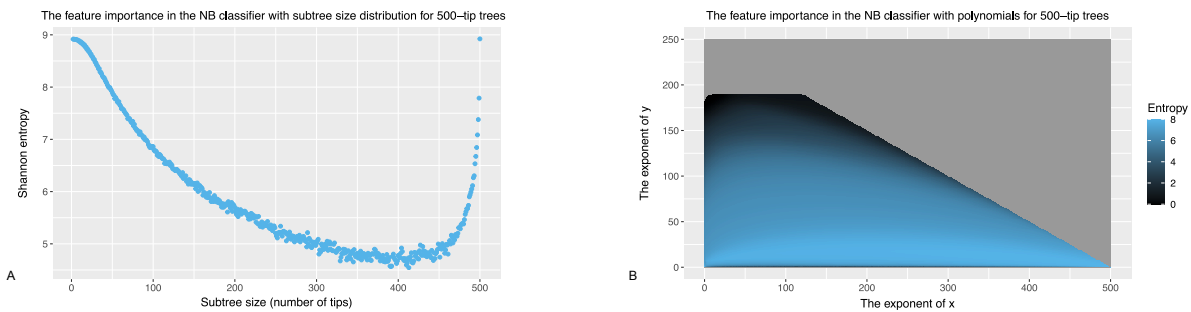
558 *WHO influenza clades* For clade 3c3.B, the 95% confidence interval of the birth  
 559 rate  $\lambda_B$  is (0.56, 0.60) and the 95% confidence interval of the death rate  $\mu$  is (0.58, 0.62).  
 560 The 95% confidence interval of  $R_0$  of the clade is (0.918, 1.013).

POLYNOMIAL PHYLOGENETIC ANALYSIS OF TREE SHAPES

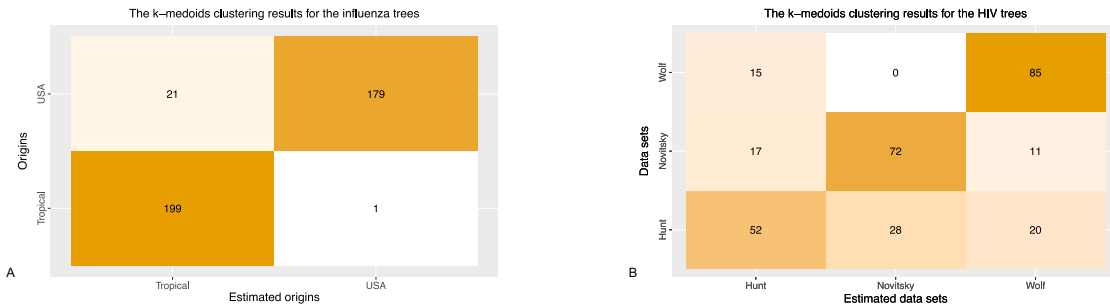
29



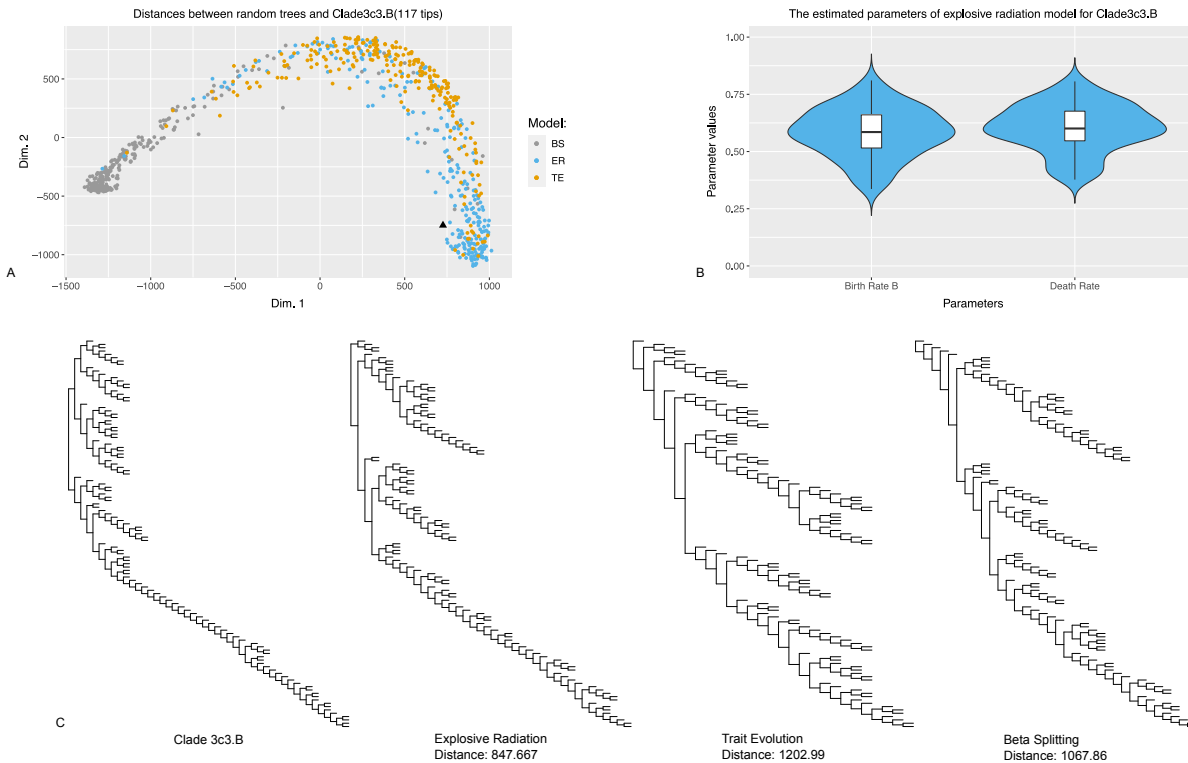
Supplementary Figure 4. A-B: the comparisons between the real parameter and the estimated parameter of the beta splitting random trees with 250 tips and 500 tips using subtree size distributions. C-D: the comparisons between the real parameters and the estimated parameters of the explosive radiation random trees with 750 tips using subtree size distributions. E-F: the results of using naive Bayes classifiers to select the model generating random trees with 500 tips using subtree size distributions.



Supplementary Figure 5. A: the feature importance (Shannon entropy) in the naive Bayes classifier used for model selection with subtree size distributions. B: the feature importance (Shannon entropy) in the naive Bayes classifier used for model selection with polynomials.



Supplementary Figure 6. A: the results of  $k$ -medoids clustering for the influenza trees using the polynomial binary differences. B: the results of  $k$ -medoids clustering for the HIV trees using the polynomial binary differences.



Supplementary Figure 7. A: the MDS plots of the polynomial distances between the random trees generated by the three different models and the clade 3c3.B. B: the distribution of the estimated parameters of the clade 3c3.B over 100 replicates. C: the plots of the clade A3 and the nearest random trees generated by the three different models.

561

REFERENCES

- 562 P. Agapow and A. Purvis. 2002. Power of eight tree shape statistics to detect nonrandom  
563 diversification: a comparison by simulation of two models of cladogenesis. *Systematic*  
564 *Biology*. 51(6):866–72.
- 565 C. Aggarwal, A. Hinneburg and D. Keim. 2001. On the surprising behavior of distance  
566 metrics in high dimensional spaces. *Proceedings of the International Conference on*  
567 *Database Theory*. 420–434.
- 568 D. Aldous. 1996. Probability distributions on cladograms. In: D. Aldous, R. Pemantle and  
569 editors, Random discrete structures. *Springer IMA Volumes in Mathematics and its*  
570 *Application*. 76:1–18.
- 571 D. Aldous. 2001. Stochastic models and descriptive statistics for phylogenetic trees, from  
572 yule to today. *Statistical Science*. 16(1):23–34.
- 573 D. Andrén and K. Markström. 2009. The bivariate Ising polynomial of a graph. *Discrete*  
574 *Appl. Math.* 157:2515–24.
- 575 T. Bedford et al.. 2015. Global circulation patterns of seasonal influenza viruses vary with  
576 antigenic drift. *Nature*. 523(7559):217–20.
- 577 L. Billera, S. Holmes and K. Vogtmann. 2001, Geometry of the space of phylogenetic trees.  
578 *Advances in Applied Mathematics*. 27(4):733–767.
- 579 M. Binet et al.. 2016. Fast and accurate branch lengths estimation for phylogenomic trees.  
580 *BMC Bioinformatics*. 17(23); doi: 10.1186/s12859-015-0821-8.
- 581 M. Blum and O. François. 2006. Which random processes describe the tree of life? A  
582 large-scale study of phylogenetic tree imbalance. *Systematic Biology*. 55(4):685–91.
- 583 P. Botti, and R. Merris. 1993. Almost all trees share a complete set of immanantal  
584 polynomials. *Journal of Graph Theory*, 17(4):467-476.



- 585 J. Brown et al.. 2010. When Trees Grow Too Long: Investigating the Causes of Highly  
586 Inaccurate Bayesian Branch-Length Estimates. *Systematic Biology*. 59(2):145–161.
- 587 S. Chaudhary and G. Gordon. 1991. Tutte polynomials for trees. *J. Graph Theory*.  
588 15:317–331.
- 589 C. Chewapreecha et al.. 2014. Dense genomic sampling identifies highways of  
590 pneumococcal recombination. *Nature Genetics* 46(3):305–309.
- 591 L. Chindelevitch et al.. 2019. Network science inspires novel tree shape statistics. *Preprint*.  
592 bioRxiv 608646; doi: <https://doi.org/10.1101/608646>.
- 593 S. Choi, S. Cha, and C. Tappert. 2010. A survey of binary similarity and distance  
594 measures. *Journal of Systemics, Cybernetics and Informatics*. 8(1):43–48.
- 595 C. Colijn and G. Plazzotta. 2018. A metric on phylogenetic tree shapes. *Systematic*  
596 *Biology*. 67:113–126.
- 597 D. Colless, 1982. Review of phylogenetics: the theory and practice of phylogenetic  
598 systematics. *Systematic Zoology*. 31(100).
- 599 A. Dayarian and B. Shraiman. 2014. How to infer relative fitness from a sample of genomic  
600 sequences. *Genetics*. 197(3):913–23.
- 601 S. Frost and E. Volz. 2013. Modelling tree shape and structure in viral phylodynamics.  
602 *Phil. Trans. R. Soc. B*. 368; doi: <http://doi.org/10.1098/rstb.2012.0208>
- 603 G. Fusco and Q. Cronk. 1995. A new method for evaluating the shape of large phylogenies.  
604 *Journal of Theoretical Biology*. 175(2):235–243.
- 605 B. Grenfell et al.. 2004. Unifying the epidemiological and evolutionary dynamics of  
606 pathogens. *Science*, 303(5656):327–332.
- 607 J. Hadfield et al.. 2018. Nextstrain: real-time tracking of pathogen evolution.  
608 *Bioinformatics*. 34(23):4121–4123.

- 609 Hartmann, K., Wong, D. and Stadler, T., 2010. Sampling trees from evolutionary models.  
610 *Systematic biology*, 59(4):465–476.
- 611 M. Hayati, P. Biller and C. Colijn. 2020. Predicting the short-term success of human  
612 influenza A variants with machine learning. *Proceedings of the Royal Society B*.  
613 287(1924):20200319.
- 614 S. Heard. 1996. Patterns in phylogenetic tree balance with variable and evolving speciation  
615 rates. *Evolution*. 50(6): 2141–2148.
- 616 G. Hunt et al.. 2013. Surveillance of transmitted HIV-1 drug resistance in 5 provinces in  
617 South Africa in 2011. *Communicable Diseases Surveillance Bulletin*. 11:122–124.
- 618 V. Jones. 1985. A polynomial invariant for knots via von Neumann algebras. *Bull. Amer.*  
619 *Math. Soc.* 12:103–111.
- 620 L. Kaufman, and P.J. Rousseeuw. 1990. Finding groups in data: An introduction to cluster  
621 analysis. New York: Wiley.
- 622 M. Kendall and C. Colijn. 2016. Mapping phylogenetic trees to reveal distinct patterns of  
623 evolution. *Molecular Biology and Evolution*. 33(10):2735–43.
- 624 M. Kendall, V. Eldholm and C. Colijn. 2018. Comparing phylogenetic trees according to  
625 tip label categories. *Preprint*. bioRxiv 251710; doi: <https://doi.org/10.1101/251710>.
- 626 M. Kirkpatrick and M. Slatkin. 1993. Searching for evolutionary patterns in the shape of a  
627 phylogenetic tree. *Evolution*. 47(4):1171–1181.
- 628 A. Lambert and T. Stadler. 2013. Birth-death models and coalescent point processes: the  
629 shape and probability of reconstructed phylogenies. *Theoretical Population Biology*.  
630 90:113–28.
- 631 E. Lewitus and H. Morlon. 2016. Characterizing and Comparing Phylogenies from their  
632 Laplacian Spectrum. *Systematic Biology*. 65(3): 495–507.

- 633 P. Liu. 2021. A tree distinguishing polynomial. *Discrete Applied Mathematics*. 288(15):1–8.
- 634 J. Losos et al.. 2013. Evolutionary biology for the 21st century. *PLoS Biology*.  
635 11(1):e1001466.
- 636 L. Maia, A Colato and J.Fontanar. 2004. Effect of selection on the topology of genealogical  
637 trees. *Journal of Theoretical Biology*. 226(3):315–20.
- 638 M. Manceau, A. Lambert and H. Morlon. 2015. Phylogenies support out-of-equilibrium  
639 models of biodiversity. *Ecology Letters*. 18(4):347–56.
- 640 F. Matsen. 2006. A geometric approach to tree shape statistics. *Systematic Biology*.  
641 55(4):652–61.
- 642 F. Matsen, and S. Evans. 2012. Ubiquity of synonymity: almost all large binary trees are  
643 not uniquely identified by their spectra or their immanantal polynomials. *Algorithms for*  
644 *Molecular Biology: AMB*. 7(1):14.
- 645 A. McKenzie and M. Steel. 2000. Distributions of cherries for two models of trees.  
646 *Mathematical Biosciences*. 164(1):81–92.
- 647 M. Monagan and B. Tuncer. 2018. Factoring multivariate polynomials with many factors  
648 and huge coefficients. *CASC*. 11077:319–34.
- 649 A. Mooers and S. Heard. 1997. Inferring evolutionary process from phylogenetic tree  
650 shape. *The Quarterly Review of Biology*. 31–54.
- 651 S. Negami and K. Ota. 1996. Polynomial invariants of graphs II. *Graphs Combin*.  
652 12:189–198.
- 653 R. Neher and T. Bedford. 2015. nextflu: Real-time tracking of seasonal influenza virus  
654 evolution in humans. *Bioinformatics*. 31(21):3546–48.
- 655 V. Novitsky et al.. 2013. Phylogenetic relatedness of circulating HIV-1C variants in  
656 Mochudi, Botswana. *PLoS One*. 8(12):e80589.

- 657 G. Plazzotta and C. Colijn. 2016. Asymptotic frequency of shapes in supercritical  
658 branching trees. *Journal of Applied Probability*. 53(4):1143–55.
- 659 M. Price, P. Dehal, and A. Arkin. 2010. Fasttree 2—approximately maximum-likelihood  
660 trees for large alignments. *PloS one*. 5(3):e9490, doi:10.1371/journal.pone.0009490.
- 661 A. Purvis et al.. 2011. The shape of mammalian phylogeny: patterns, processes and scales.  
662 *Philosophical Transactions of the Royal Society B*. 366(1577):2462–77.
- 663 O. Pybus and P. Harvey. 2000. Testing macro-evolutionary models using incomplete  
664 molecular phylogenies. *Proc. R. Soc. Lond. B*. 267:2267–2272.
- 665 I. Rish. 2001. An empirical study of the naive Bayes classifier. *Proceedings of the IJCAI-01*  
666 *Workshop on Empirical Methods in Artificial Intelligence* 41–46.
- 667 D. Robinson and L. Foulds. 1981. Comparison of phylogenetic trees. *Mathematical*  
668 *Biosciences*. 53(1-2):131–47.
- 669 C. Russell et al.. 2008. The global circulation of seasonal influenza a (H3N2) viruses.  
670 *Science*. 320:340–46 .
- 671 M. Sackin. 1972. “Good” and “bad” phenograms. *Systematic Zoology*. 21(2):225–26.
- 672 R. Safavian, and D. Landgrebe. 1991. A survey of decision tree classifier methodology.  
673 *IEEE Transactions on Systems, Man, and Cybernetics*. 21(3):660–74.
- 674 E. Saulnier , O. Gascuel, and S. Alizon. 2017. Inferring epidemiological parameters from  
675 phylogenies using regression-ABC: A comparative study. *PLOS Computational Biology*.  
676 13(3):e1005416.
- 677 A. Stamatakis. 2014. RAxML version 8: a tool for phylogenetic analysis and post-analysis  
678 of large phylogenies. *Bioinformatics*. 30(9):1312–13.
- 679 M. Steel and A. McKenzie. 2001. Properties of phylogenetic trees generated by Yule-type  
680 speciation models. *Mathematical biosciences*. 170(1):91–112.

- 681 M. Stich and S. Manrubia. 2009. Topological properties of phylogenetic trees in  
682 evolutionary models. *The European Physical Journal B*. 70(4):583–92.
- 683 T. To et al.. 2016. Fast Dating Using Least-Squares Criteria and Algorithms. *Systematic*  
684 *Biology* 65(1):82–97.
- 685 W. Tutte. 1954. A contribution to the theory of chromatic polynomials. *Can. J. Math.*  
686 6:80–91.
- 687 L. van der Maaten and G. Hinton. 2008. Visualizing High-Dimensional Data Using t-SNE.  
688 *Journal of Machine Learning Research*. 9(11):2579–2605.
- 689 E. Volz, K. Koelle and T. Bedford. 2013. Viral phylodynamics. *PLoS Computational Biology*.  
690 9(3):e1002947.
- 691 E. Wolf et al.. 2017. Phylogenetic evidence of HIV-1 transmission between adult and  
692 adolescent men who have sex with men. *AIDS Research and Human Retroviruses*.  
693 33:318–22.
- 694 T. Wu and K. Choi. 2016. On joint subtree distributions under two evolutionary models.  
695 *Theoretical Population Biology*. 108:13–23.