

Gene Transfer-based Phylogenetics: Analytical Expressions and Additivity via Birth–Death Theory

GUY KATRIEL¹, UDI MAHANAYMI², SHELLY BREZNER², NOOR KEZEL³ CHRISTOPH KOUTSCHAN⁴,
DORON ZEILBERGER⁵, MIKE STEEL⁶, AND SAGI SNIR^{2*}

¹*Department of Mathematics, Ort Braude, Israel*

²*Department of Evolutionary and Environmental Biology, University of Haifa, Israel*

³*Department of Mathematics, University of Haifa, Israel*

⁴*RICAM, Austrian Academy of Sciences, Linz, Austria*

⁵*Department of Mathematics, Rutgers University, USA*

⁶*School of Mathematics and Statistics, University of Canterbury, NZ.*

**Corresponding author: ssagi@research.haifa.ac.il*

ABSTRACT

1 The genomic era has opened up vast opportunities in molecular systematics, one of which
2 is deciphering the evolutionary history in fine detail. Under this mass of data, analyzing
3 the point mutations of standard markers is often too crude and slow for fine-scale
4 phylogenetics. Nevertheless, genome dynamics (GD) events provide alternative, often
5 richer information. The *synteny index* (SI) between a pair of genomes combines gene order
6 and gene content information, allowing the comparison of genomes of unequal gene
7 content, together with order considerations of their common genes. Recently, genome
8 dynamics has been modelled as a continuous-time Markov process, and gene distance in
9 the genome as a birth–death–immigration process. Nevertheless, due to complexities
10 arising in this setting, no precise and provably consistent estimators could be derived,
11 resulting in heuristic solutions.

12 Here, we extend this modelling approach by using techniques from birth–death
13 theory to derive explicit expressions of the system’s probabilistic dynamics in the form of
14 rational functions of the model parameters. This, in turn, allows us to infer analytically
15 accurate distances between organisms based on their SI. Subsequently, we establish
16 additivity of this estimated evolutionary distance (a desirable property yielding
17 phylogenetic consistency).

18 Applying the new measure in simulation studies shows that it provides accurate
19 results in realistic settings and even under model extensions such as gene gain/loss or over
20 a tree structure. In the real-data realm, we applied the new formulation to unique data
21 structure that we constructed - the ordered orthology DB - based on a new version of the
22 EggNOG database, to construct a tree with more than 4.5K taxa. To the best of our
23 knowledge, this is the largest gene-order-based tree constructed and it overcomes
24 shortcomings found in previous approaches. Constructing a GD-based tree allows to
25 confirm and contrast findings based on other phylogenetic approaches, as we show.

26 *Key words:* **keywords:** Genome Dynamics, Prokaryotic Phylogenetics, Statistical
27 Consistency, Synteny Index

INTRODUCTION

The genomic era has reached the point where tasks that seemed imaginary only a decade ago are now within reach. Among these tasks is the inference of the evolutionary history for tens of thousands of species, sometimes of very close origin. Such a history is depicted in a tree structure and is called a *phylogeny* or a *phylogenetic tree*. The leaves of that tree correspond to contemporary extant species, internal nodes correspond to ancestral species, and the tree's edges (or branches) between nodes correspond to evolutionary relationships. Despite the impressive advances in the extraction of molecular data, and of ever-increasing quality, finding a phylogenetic tree which accounts for the data in a satisfactory way is still a major challenge that requires reliable approaches for inferring the true evolutionary relationships between the species under study.

Statistical modelling in which the tree is a parameter of some assumed model is nowadays considered the method of choice for phylogenetic inference. Under this framework, vast efforts have been made, first to model data accurately, and then to draw inferences efficiently from the given data. One such approach is *maximum likelihood* (Felsenstein, 1978, 1981; Hasegawa et al., 1991; Yang, 1996), where the model (tree) selected is the one maximising the probability of observing the given data.

Standard phylogenetics, whether parsimony- or likelihood-based, analyses one or a few ubiquitous genes residing in all species under study, and uses the differences between respective gene copies i.e., orthologues, in order to infer evolution history. These genes are typically highly conserved by definition and therefore advantageous in certain settings such as very rapid viral evolution (Pybus and Rambaut, 2009) or long evolutionary distance where finer markers saturate (Ciccarelli et al., 2006). However, for the task of distinguishing the shallow branches of the prokaryotic tree, these genes often fail to provide a strong enough signal (Sevillya and Snir, 2019; Martinez-Gutierrez and Aylward, 2021; Rajendhran and Gunasekaran, 2011). In contrast, genome dynamics events (GDE's) are larger scale events compared to single nucleotide mutations, in which a complete gene or a sequence of genes, are involved. One such event is horizontal gene transfer (HGT), a mechanism by which organisms transfer genetic material to contemporaneous organisms rather than via vertical inheritance (Doolittle, 1999; Koonin et al., 2001; Ochman et al., 2000). Among prokaryotes, GDE's in the form of HGT and gene loss seem to provide far richer information, as indicated in (Puigbò et al., 2014): "The rates of genome change are remarkably high, typically tens of thousands of GDEs per nucleotide substitution per site, or tens to hundreds of GDEs per substitution per gene", and see also in e.g. (Schnknecht et al., 2014; Pang and Lercher, 2019; Koonin et al., 2021). The latter fact calls for GD-based phylogenetic approaches, in particular when handling prokaryotes sharing a close origin. GD-based phylogenetics is mainly divided into gene-order-based and gene-content-based techniques. With the gene-order-based approach (Sankoff, 1992; Hannenhalli and Pevzner, 1999; Yancopoulos et al., 2005), two genomes are considered as permutations of the gene set, and distance is defined as the minimal number of operations needed to transform one genome to the other. In the other, the gene-content-based approach (Snel et al., 1999; Tekaiia and Dujon, 1999; Fitz Gibbon and House, 1999) gene order is entirely ignored, and similarity is defined as the size of the set of shared genes. A statistical framework has been devised for part of both these models, the order- and content-based (Serdoz et al., 2017; Wang and Warnow, 2001; Biller et al., 2015; Sankoff and Nadeau, 1996; Lin et al., 2013; Zhao et al., 2021). The Jump operation studied here is accounted for by some these gene-order models; however to the best of our knowledge, a stochastic, rate-dependent framework accounting for HGT, has not been suggested.

Another and related line of works relies on gene tree amalgamation into a unified species tree, a task referred to as a *supertree construction* (Strimmer and Moulton, 2000; Bininda-Emonds, 2004; Baum, 1992; Ragan, 1992). Some of these works are

80 likelihood-based where GD events are accounted for by the differences in gene tree
81 topologies (Morel et al., 2022, 2020). Nevertheless, these latter works ignore gene order,
82 and in particular gene order likelihood.

83 Thus, devising a genuine evolutionary model, along with an estimator of the model
84 parameters from observed data only, and an efficient inference method of this estimator,
85 remains a challenging task.

86 A related task in this field is the *reconciliation* between a gene tree and the species
87 tree. In this setting, a sequence of events acting on the species tree and yielding the given
88 gene tree is sought. These events may contain events other than HGT which are commonly
89 denoted *duplication*, *transfer* and *loss* (DTL) (Bansal et al., 2018; Stolzer et al., 2012).
90 These works contain both parsimony-based approaches such as (Nakhleh et al., 2005;
91 Doyon et al., 2010), and model/likelihood-based approaches (Szöllősi et al., 2013;
92 Sjöstrand et al., 2014). Although it deals with the same objects, and the same events, as
93 the tree is already given, the goal there is not tree reconstruction, and in particular not
94 reconstruction based on gene order between multiple genes, as in our case here.

95 The *synteny index* (SI) (Shifman et al., 2013; Adato et al., 2015) was suggested as
96 an alternative measure to the parsimony/statistical phylogenetics approaches mentioned
97 above, allowing unequal gene content on one hand while accounting for the order among
98 the shared genes. Here, the locality of a gene in the form of a “neighbourhood” is
99 considered and compared with other genomes. Similarity between genomes is attained by
100 averaging this local quantity over all the shared genes.

101 Aiming at a rigorous delineation of SI, in a recent paper (Sevillya et al., 2019), we
102 defined an underlying simplistic model, the *Jump model*, to model genome dynamics,
103 primarily HGT. Under this model, every gene stochastically (at some constant rate)
104 “jumps” to a random location at the chromosome. Consequently distance between two
105 genes along the genome, i.e., the number of genes separating between these two genes, can
106 be described as a (critical) birth–death–immigration process. The setting poses intrinsic
107 hurdles such as overlapping neighbourhoods, non-stationarity, confounding factors, and
108 more. Consequently, precise quantities could not be obtained in this earlier work and were
109 calculated heuristically. The Jump model consists of a Jump operation embedded within a
110 stochastic framework. While the basic Jump operation is subsumed in some of the
111 gene-order models mentioned above, to the best of our knowledge, no complete
112 time-dependent framework accounting for HGT, has been suggested. (Dalevi and Eriksen,
113 2008) defines the single gene transposition model that is equivalent to a Jump, and
114 expected distances are derived by a function of the number of breakpoints, however, the
115 meaning of a model there is a type of operation as opposed to a stochastic, rate-dependent
116 model considered here.

117 In this work, we take the Jump model and the SI a significant step further by
118 deriving exact and invertible analytical expressions (Theorems 2, 3, and 4) that allow for
119 the evolutionary distance between species to be inferred from the (averaged) SI values
120 under the Jump model. By an earlier result (Theorem 1), this implies that the difference
121 between these estimates of evolutionary distance and the true evolutionary distances
122 converges to zero as the number of genes becomes large. Our results rely on techniques
123 from the theory of birth-death processes. On the experimental side, we first show that the
124 new expressions provide accurate reconstructions, even for real-life problem sizes although
125 the theoretical underlying model on which these expressions are based assumes infinitely
126 long genomes. We note here that, for the sake of rigorous analysis, the pure theoretical
127 Jump model consists of only the Jump operation, and hence implies comparisons between
128 equal content genomes - genomes over the same gene set. Such a model can can
129 accommodate other, however restricted, evolutionary scenarios including gene gain/loss
130 events, as we show in the Methods part. Nevertheless, to allow for a broader range of
131 scenarios resulting in genomes with unequal content, as is the case in real data, we have

132 devised two heuristics - the *union* and the *intersection* gene set approaches described in
 133 the experimental section. Using the Jump model under these heuristics to simulated data
 134 including both jumps and gene gain and loss events (i.e., unequal gene content), shows
 135 robustness to such more diverse regimes.

136 For real data, we created a new database of ordered orthology groups, based on the
 137 EggNOG (Huerta-Cepas et al., 2018) orthology database, encompassing over 4445
 138 organisms spanning the entire prokaryotic phylogenetic spectrum. Applying the new
 139 measure to this database, produces a tree with very high agreement with the NCBI
 140 taxonomy (Federhen, 2011; Schoch et al., 2020). To the best of our knowledge, this is the
 141 largest genome-dynamics-based tree. In comparison with other SI-based trees, it is evident
 142 that the new technique reconstructs significantly more realistic distances, attesting to its
 143 capability as a distance measure in various other applications of genome dynamics (Che
 144 et al., 2006; Rogozin et al., 2002). Moreover, contrasting the SI-based trees with the NCBI
 145 taxonomy, suggests several incongruences that may be of independent, intrinsic interest.

146 **Comment:** For the sake of readability, technical proofs have been moved to the
 147 Appendix. Additionally, as both the theoretical and the experimental parts are technically
 148 involved, we provide in the Supplementary Text brief self-contained background to the
 149 theoretical material employed, as well as further details for the experimental parts. Finally,
 150 for the sake of reconstructability, we provide in the Supplementary Material data produced
 151 during this research. Supplementary Text and Material are found in the DRYAD link.

152 MATERIALS AND METHODS

153 We start by defining a restricted model – *the Jump model* – which can be regarded
 154 as a transfer between genomes over the same gene set (*equal content*). Biologically, the
 155 *Jump* operation, in which a gene moves to another location, can account for several GDE’s,
 156 such as a gene duplication and a subsequent loss, a gene loss in which a gene jumps outside
 157 of the genome, a gene gain when the Jump is from an alien genome, or both (gain and
 158 subsequent loss, or vice versa) as discussed in e.g. (Liu et al., 2004).

159 *The Jump Model:* Let $\mathcal{G}^n = (g_1, g_2, \dots, g_n)$ be a sequence of n ‘genes’, and
 160 henceforth we remove the superscript n as it holds throughout. In the analysis, we will
 161 assume that n is large enough to allow us to ignore the tips of \mathcal{G} (or, equivalently, \mathcal{G} is
 162 cyclic and there are no tips). We now introduce a stochastic process operating on \mathcal{G} .
 163 Consider the following continuous-time Markovian process $\mathcal{G}(t), t \geq 0$ on the state space of
 164 all $n!$ permutations of g_1, g_2, \dots, g_n . Each gene g_i is independently subject to a Poisson
 165 process transfer event (at a constant rate λ) in which g_i is moved to a different position in
 166 the sequence, with each of the possible $n - 1$ positions (between consecutive genes that are
 167 different from g_i , or at the start or end of the sequence) and with this target location for
 168 the transfer selected uniformly at random from these $n - 1$ possibilities.

169 For example, if $\mathcal{G}(t) = (g_1, g_2, g_3, g_4, g_5)$, then g_4 might transfer to be inserted
 170 between g_1 and g_2 to give the sequence $\mathcal{G}(t + \delta) = (g_1, g_4, g_2, g_3, g_5)$. The other sequences
 171 that could arise by a single transfer of g_4 are $(g_4, g_1, g_2, g_3, g_5)$, $(g_1, g_2, g_4, g_3, g_5)$, and
 172 $(g_1, g_2, g_3, g_5, g_4)$.

173 Since the model assumes a Poisson process, the probability that g_i is transferred to
 174 a different position between times t and $t + \delta$ is $\lambda\delta + o(\delta)$, where the $o(\delta)$ term accounts for
 175 the possibility of more than one transfer occurring in the time period δ (this possibility has
 176 probability of order δ^2 and so is asymptotically negligible compared with terms of order δ
 177 as $\delta \rightarrow 0$). Moreover, a single transfer event always results in a different sequence.

178 *The Synteny Index:* Let k be any constant positive integer (note that it may be
 179 possible to allow k to grow slowly with n , but we will not explore such an extension here).
 180 For $j \in k + 1, \dots, n - k$, the $2k$ -neighbourhood of gene g_j in a genome \mathcal{G} , $N_{2k}(g_j, \mathcal{G})$ is the
 181 set of $2k$ genes (different from g_j) that have a distance of at most k from g_j in \mathcal{G} . We also
 182 define $SI_j(t)$ as the relative intersection size between $N_{2k}(g_j, \mathcal{G}(0))$ and $N_{2k}(g_j, \mathcal{G}(t))$, or
 183 formally, $SI_j(t) = \frac{1}{2k} |N_{2k}(g_j, \mathcal{G}(0)) \cap N_{2k}(g_j, \mathcal{G}(t))|$ (this is also called *the Jaccard index*
 184 between the two neighbourhoods (Jaccard, 1901)).

185 Let $\overline{SI}(\mathcal{G}(0), \mathcal{G}(t))$ be the average of these $SI_j(t)$ values over all j between $k + 1$ and
 186 $n - k$. That is,

$$\overline{SI}(\mathcal{G}(0), \mathcal{G}(t)) = \frac{1}{n - 2k} \sum_{j=k+1}^{n-k} SI_j(t). \quad (1)$$

187 Subsequently, when time t does not matter, we simply use \overline{SI} or simply SI where it
 188 is clear from the context.

189 *Phylogenetic Trees and Distances:* For a set of species (denoted *taxa*) \mathcal{X} , a
 190 phylogenetic \mathcal{X} -tree T is a tree $T = (V, E)$ for which there is a one-to-one correspondence
 191 between \mathcal{X} and the set $\mathcal{L}(T)$ of leaves of T . A tree T is *weighted* if there is a weight (or
 192 length) function associating non-negative weights (lengths) to the edges of T . In this
 193 paper, we will use the term length, as it corresponds to the number of events or the time
 194 span. Edge lengths are naturally extended to *paths*, where the path length is the sum of
 195 edge lengths along the path. Assume a model M operating on T by generating events,
 196 starting from the root down to the leaves, where edge lengths serve as expected number of
 197 events generated by M . The notion of additivity is classical in phylogenetics (Buneman,
 198 1971; Semple and Steel, 2003). We briefly specialise it to our case. Let D be some *pairwise*
 199 *marker* or *a distance measure* between the outcome of M on the leaves of T . Then D is
 200 said to be *additive on the model M* if D can be transformed (or *corrected*) by applying a
 201 fixed function f to D (only), such that the corrected distance converges to the expected
 202 number of events under the model M , as the amount of data (e.g. sequence length, or in
 203 our case n , the number of genes), becomes large.

204 *Gene Neighbourhood as a Markov Chain*

205 We now introduce a local random process, induced by the Jump model defined
 206 above that operates on the entire genome level. This local model will play a key role in the
 207 analysis of the random variable $\overline{SI}(\mathcal{G}(0), \mathcal{G}(t))$. Consider the location of a gene g_i , that
 208 does not jump during time period t , with respect to another gene $g_{i'}$. Without loss of
 209 generality assume $i > i'$ and let $j = i - i'$. Now, there are j ‘slots’ between $g_{i'}$ and g_i into
 210 which a third jumping gene g_k can be inserted, but only $j - 1$ genes in that interval can
 211 jump out. Obviously, a jump into that interval moves $g_{i'}$ one position away from g_i , and a
 212 jump from that interval, moves $g_{i'}$ closer to g_i . This local model, describing the distance
 213 between $g_{i'}$ and g_i , can be described by a continuous-time random walk on the state space
 214 $1, 2, 3, \dots$ with transitions from j to $j + 1$ at rate $j\lambda$ (for all $j \geq 1$) and from j to $j - 1$ at
 215 rate $(j - 1)\lambda$ (for all $j \geq 2$), with all other transition rates being 0. This, the distance
 216 between $g_{i'}$ and g_i , is thus a (generalised linear) birth–death process, illustrated in Fig. 1.
 217 We note though that this is not an independent model, occurring between individual
 218 genes, separated from the Jump model operating at the genome scale. Rather, the
 219 birth–death process modeling the distance between individual genes, is induced by that
 220 larger scale model of the Jump.

More formally, we will let X_t denote the random variable that describes the number of slots between two genes under this process described above. Then X_t is a continuous-time random walk on state space $1, 2, 3, \dots$, with an arbitrary initial condition X_0 and transition probabilities of X_t defined as follows:

$$\mathbb{P}(X_{t+\delta} = j + 1 | X_t = j) = j\lambda\delta + o(\delta), \quad j \geq 1, \quad (2)$$

$$\mathbb{P}(X_{t+\delta} = j - 1 | X_t = j) = (j - 1)\lambda\delta + o(\delta), \quad j \geq 1. \quad (3)$$

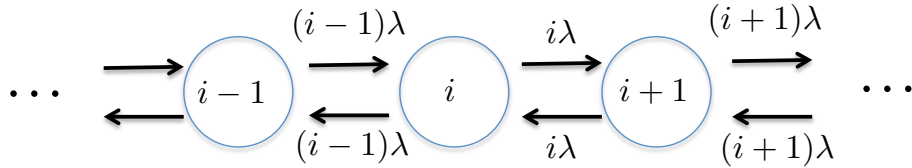


Fig. 1. Transitions for the process X_t

221 The process X_t is slightly different from the much-studied critical linear birth–death
 222 process, for which the rate of birth and death from state j are both equal to j (here the
 223 rate of birth is j but the rate of death is $j - 1$), and for which 0 is an absorbing state (here
 224 there are no absorbing states). However, this stochastic process is essentially a translation
 225 of a critical linear birth–death process with immigration rate equal to the birth–death
 226 rate λ . This connection is key to the analysis of the divergence times that we establish
 227 below.

228 RESULTS

229 *Explicit Expressions for the Divergence Time*

230 We now present the main theoretical contribution of this work, which is an
 231 analytical expression of divergence times. We first recall a result of (Sevillya et al., 2019)
 232 that links SI and the transition probabilities of the birth–death process X_t . This raises the
 233 need to obtain explicit expressions for these probabilities, which we do by making use of
 234 known results from the theory of birth–death processes. This theory also allows us to give
 235 a proof of the monotonicity of the SI as a function of time (in the limit of large n), a result
 236 that is crucial in order to ensure that we can use our explicit expressions to solve the
 237 divergence time in terms of the SI.

Let $p_{i,j}(t)$ be the transition probability for X_t to be at state j , given that at time 0 it was at state i :

$$p_{i,j}(t) = \mathbb{P}(X_t = j | X_0 = i), \quad i, j \geq 1.$$

238 We denote the conditional probability that $X_t \in [k]$ given that $X_0 = i$ by:

$$q_{i,k}(t) = \sum_{j=1}^k p_{i,j}(t). \quad (4)$$

239 Next, let

$$q_k(t) := \frac{1}{k} \sum_{i=1}^k q_{i,k}(t) = \frac{1}{k} \sum_{i=1}^k \sum_{j=1}^k p_{i,j}(t). \quad (5)$$

240 The quantity $q_k(t)$ is the probability that for a gene at an initial state i (i.e., at
 241 distance from a reference gene) chosen uniformly at random between 1 and k , the process
 242 X_* is still between 1 and k at time t . In (Sevillya et al., 2019) we proved the following
 243 result:

THEOREM 1 For any given value of t , as $n \rightarrow \infty$:

$$\overline{SI}(\mathcal{G}(0), \mathcal{G}(t)) \xrightarrow{p} \exp(-2\lambda t)q_k(t),$$

244 where \xrightarrow{p} denotes convergence in probability.

245 In the following we assume, without loss of generality, that $\lambda = 1$ (this is simply
 246 rescaling time). The functions $p_{i,j}(t)$ can be expressed as solutions of an infinite system of
 247 ordinary differential equations (Sevillya et al., 2019) (the Kolmogorov forward equations
 248 corresponding to the birth–death process), and these differential equations may be used to
 249 numerically approximate $p_{i,j}(t)$ and therefore the key quantity $q_k(t)$. However, in the
 250 present paper we will derive *explicit* algebraic expressions for $p_{i,j}(t)$ and thus $q_k(t)$. It
 251 thereby becomes possible to use Theorem 1 to solve for the divergence time t in terms of
 252 the SI.

253 *Explicit expressions for $p_{i,j}(t)$*

THEOREM 2

$$p_{i,j}(t) = \frac{1}{(t+1)^{i+j-1}} \cdot \sum_{\ell=1}^{\min(i,j)} \frac{(i+j-\ell-1)!}{(i-\ell)!(j-\ell)!(\ell-1)!} (1-t^2)^{\ell-1} t^{i+j-2\ell}. \quad (6)$$

254 This result follows from some general results for birth–death processes
 255 (see (Anderson, 2012) for more details). The full proof is given in the Appendix.

256 *Explicit Expression for $q_k(t)$*

257 As stated above, Theorem 1 (originally from (Sevillya et al., 2019)) gives an
 258 expression for the SI value between two genomes, $\mathcal{G}(0)$ and $\mathcal{G}(t)$. Nevertheless, in that
 259 paper, we could not derive an expression only in terms of the number of events that
 260 occurred during time t (or, alternatively, in a path along the tree of length λt “separating”
 261 genomes \mathcal{G}_i and \mathcal{G}_j) as we could not arrive at an explicit expression for q_k . Now that we
 262 have obtained explicit expression for $p_{i,j}(t)$ in Theorem 2 we can explicitly describe q_k as
 263 follows.

THEOREM 3

$$q_k(t) = \frac{1}{k} \sum_{\ell=0}^{k-1} \sum_{i=0}^{k-\ell-1} \sum_{j=0}^{k-\ell-1} \frac{(i+j+\ell)!}{i!j!\ell!} t^{i+j} (t+1)^{-i-j-2\ell-1} (1-t^2)^\ell. \quad (7)$$

264 The proof is brought in the Appendix.

To give Theorem 3 an actual expression, we provide a few instances of the above
 formula:

$$q_2(t) = \frac{2t^2 + 2t + 1}{(t+1)^3}$$

$$\begin{aligned}
q_3(t) &= \frac{3t^4 + 6t^3 + 8t^2 + 4t + 1}{(t+1)^5} \\
q_4(t) &= \frac{4t^6 + 12t^5 + 26t^4 + 26t^3 + 18t^2 + 6t + 1}{(t+1)^7} \\
q_5(t) &= \frac{5t^8 + 20t^7 + 60t^6 + 90t^5 + 102t^4 + 68t^3 + 32t^2 + 8t + 1}{(t+1)^9}.
\end{aligned}$$

265 In the supplementary text we provide the actual function for $q_{10}(t)$ that was used in
266 the real data analysis.

267 *Monotonicity of the SI Measure*

268 Recall that we assumed, without loss of generality, that $\lambda = 1$, and so our goal now
269 is to prove the monotonicity of the function, $h_k(t) = e^{-2t}q_k(t)$ and thus (by Theorem 1) the
270 SI measure itself, in the limit of large n . In fact we will prove that $q_k(t)$ itself is monotone
271 decreasing, which obviously implies that $h_k(t)$ is also monotone decreasing.

272 **THEOREM 4** The function $q_k(t)$ is monotone decreasing on $[0, \infty)$.

273 The fact that $h_k(t) = \exp(-2t)q_k(t)$ is strictly monotone decreasing with t implies
274 that $h_k(t)$ is a one-to-one-function (or *injective*) and hence the inverse function h_k^{-1} is
275 well-defined. This allows us to use Theorem 1 to reconstruct the expected number of
276 events (i.e. jump, or equivalently the time t) separating two sequences of n genes (where n
277 is large) given their pairwise SI value. By applying h_k^{-1} to the \overline{SI} for these two gene
278 sequences, the estimated number of events is obtained. By Theorem 3, we have an explicit
279 value for $h_k(t)$, so the value $h_k^{-1}(\overline{SI})$ can be calculated by numerically solving a simple
280 equation.

281 Since the expected number of events is additive on the tree, that is, the sum of
282 events along the tree edges separating two leaves equals the number of events occurred
283 between these leaves, we can conclude the following corollary:

284 **Corollary 1** Assume a genome with n genes at the root of an underlying tree T , is evolving
285 according to the Jump model defined above. Then, as $n \rightarrow \infty$, the number of events per
286 gene between the leaves of T can be reconstructed in a statistically consistent way from the
287 \overline{SI} values between the genomes at the leaves of T , by applying the transformation h_k^{-1} .

288 As a fully resolved (i.e. binary) tree T can be uniquely and consistently
289 reconstructed from its pairwise distances by applying the distance-based reconstruction
290 method, e.g. the Neighbour-Joining (NJ) algorithm, we obtain the following:

291 **Corollary 2** Assume a genome is evolving on a tree binary tree T as in Corollary 1. Then T
292 can be reconstructed in a statistically consistent way (as n grows) by transforming the \overline{SI}
293 values.

294 *Experimental Results*

295 In this section, we describe the experiments we conducted to demonstrate the
296 applicability of the theoretical results described above. We begin with simulation results
297 based on the Jump model and then move to an analysis of real genomic data.

298 *Simulation Results, Single Edge, the Pure Jump Model:* We simulated the Jump
 299 model over a single edge of length t , i.e. from $\mathcal{G}(0)$ to $\mathcal{G}(t)$, and for various values of the
 300 number n of genes. We set $k = 3$ (i.e. a neighbourhood of $2k = 6$), the rate was fixed at
 301 $\lambda = 1$ and time t varied over the interval $[0, 0.5]$. This has yielded a Jump probability that
 302 was applied to every gene in the initial (parent, $t = 0$) genome. For each value of t , the SI
 303 between the parent and the child genome was computed. The top part of Fig. 2 displays
 304 the value $e^{-2t}SI(t)$ (recall that $\lambda = 1$ and hence vanishes at the exponent) for each of 10
 305 simulations, and the function $q_3(t)$ which is the limit to which $e^{-2t}SI(t)$ converges as
 306 $n \rightarrow \infty$. As can be seen, although there is some variability due to randomness, this
 307 variability decreases as n increases, and the agreement with the limiting curve $q_3(t)$ is clear.

308 In a related experiment, we checked how well the value $SI(t)$, computed using the
 309 simulated data, can be used to estimate the time t . For each value of t , we compute $SI(t)$
 310 from the simulated data, and use this to estimate t by numerically solving the following
 311 equation:

$$e^{-2\hat{t}}q_3(\hat{t}) = SI(t), \quad (8)$$

312 In the lower part of Fig. 2 the true value of t is compared with the estimated values \hat{t} for
 313 10 simulations.

314 We note that the relevant values of λt as found in (Sevillya et al., 2019) are around
 315 0.4 for distances within the phylogenetic rank of genus. We see that the error is almost
 316 insignificant even for realistic genome sizes, as we have here.

317 *Simulation Results, Single Edge, Adding Gene Gain and Loss:* Next, we extended
 318 the pure Jump model to include gain/loss events; still as above over a single edge: for each
 319 jump event, we also generate a gain/loss event with probability p , with equal probability
 320 for a gain or a loss, so that the expected genome length is fixed. However, here we face the
 321 problem that the gene content of pairs of genomes are not identical, a fact which needs to
 322 be accounted for when computing the SI of two genomes. We have devised two different
 323 approaches to computing the SI of two genomes with non-identical gene content.

324 In the first approach (I), *the union gene set approach*, we simply replace the sum in
 325 (1) by a sum only over the genes that are common to both genomes; however, the k
 326 neighborhoods whose intersection is used to define the quantities $SI_j(t)$ include also genes
 327 which are present in only one of the two genomes. We again used Eq. (7) to infer the
 328 distances. The results for the gain/loss probabilities $p = 0.1$ and $p = 0.2$ are shown in
 329 Fig. 3.

330 In a second approach (II), *the intersection gene set approach*, for computing the SI
 331 between two genomes, we first excise all the genes which are not common to both genomes
 332 from each of the genomes. This produces a pair of genomes of the same size (i.e., the size
 333 of the intersection of the genome contents of the original genomes), and the two genomes
 334 now have identical gene content, so that their SI can be computed in the standard way.
 335 The results for the gain/loss probabilities $p = 0.1$ and $p = 0.2$ are shown in Fig. 4.

336 Comparing Fig. 3 and Fig. 4, it is clear that the second approach for computing the
 337 SI of genomes with non-identical gene content is the more appropriate one for enabling
 338 accurate estimation in the presence of gene gain and loss; indeed in Fig. 3 we see that
 339 there is a systematic bias in the estimators, which is not present in Fig. 4. Of course the
 340 procedure of excising all genes which are not common to both genomes, performed in the
 341 second approach, entails some loss of information, which is responsible for the larger
 342 variance of the estimators as seen in Fig. 4 compared to the case $n = 2000$ shown at Fig. 2.

343 *Simulation Results over Tree Structure:* Our last extension in this part is from a
 344 single edge (as reported above) to a tree structure. We describe this briefly here (a detailed

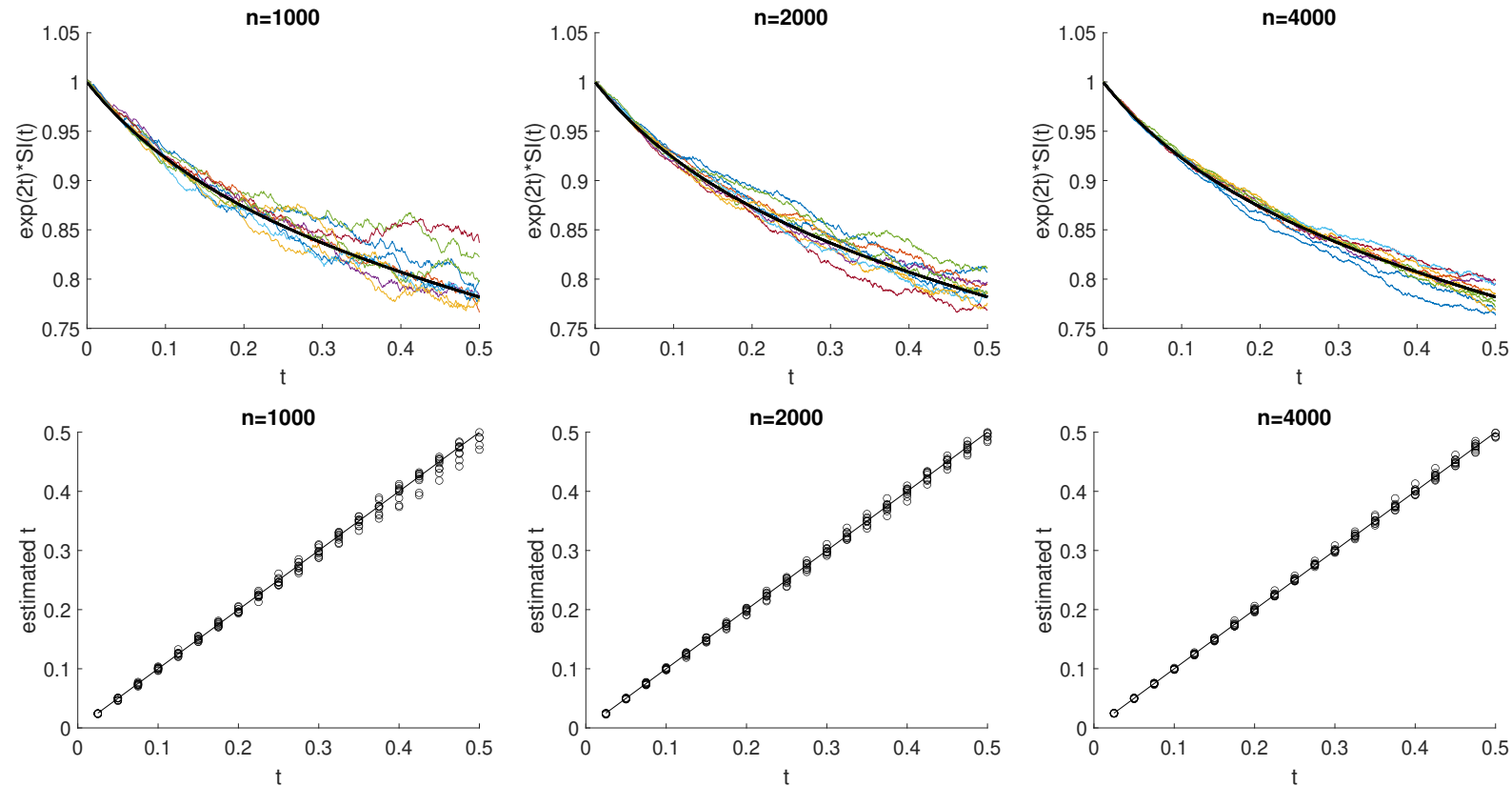


Fig. 2. **Simulation Results, Single Edge, the pure jump model:** Genome sizes $n = 1000, 2000, 4000$. Top: Comparison of the curves $e^{\lambda t}SI(t)$ computed using simulation with the limiting curve $q_3(t)$, Bottom: Estimated vs effective t .

345 description of the procedure is provided in the supplementary text). We first draw a
 346 random tree with edge lengths distributed exponentially with mean l . A genome evolves
 347 recursively down this tree starting at the root with the identity permutation, via GD
 348 events such as Jump, gain, and loss, identically as the *single-edge* experiments described
 349 above. For each edge, the edge length e_l is the expected number of GD events. At the end
 350 of this procedure, we have genomes at the leaves over ordered subsets of the initial set at
 351 the root. We applied the intersection gene set approach (Approach II above) to cope with
 352 presence of gene gain/loss, in order to reconstruct a tree. The reconstructed tree was
 353 compared to the original model tree, using the Robinson–Foulds (RF) symmetric
 354 difference (Robinson and Foulds, 1981). The results, in terms of normalized error rate
 355 (incorrect edges) versus average edge length, are shown in Fig. 5 for a tree over 26 leaves.
 356 As can be seen, for small values of edge lengths, reconstruction quality is fairly high,
 357 almost perfect. Nevertheless tree distance rises (i.e. reconstruction quality falls) sharply
 358 initially and slowly levels off towards the value of 1. Note though that even at average jump
 359 rate of one, we still observe a reconstruction of 0.5 meaning half of the edges are correct.

360 *Real Data Results:* Here we report the real data results obtained using the new
 361 technique. Because of space limitations, and for the sake of reconstructability, fuller details
 362 and data are provided in the supplementary text and material respectively. We applied our
 363 method to real genomic data consisting of 4445 prokaryotes taken from the orthology data

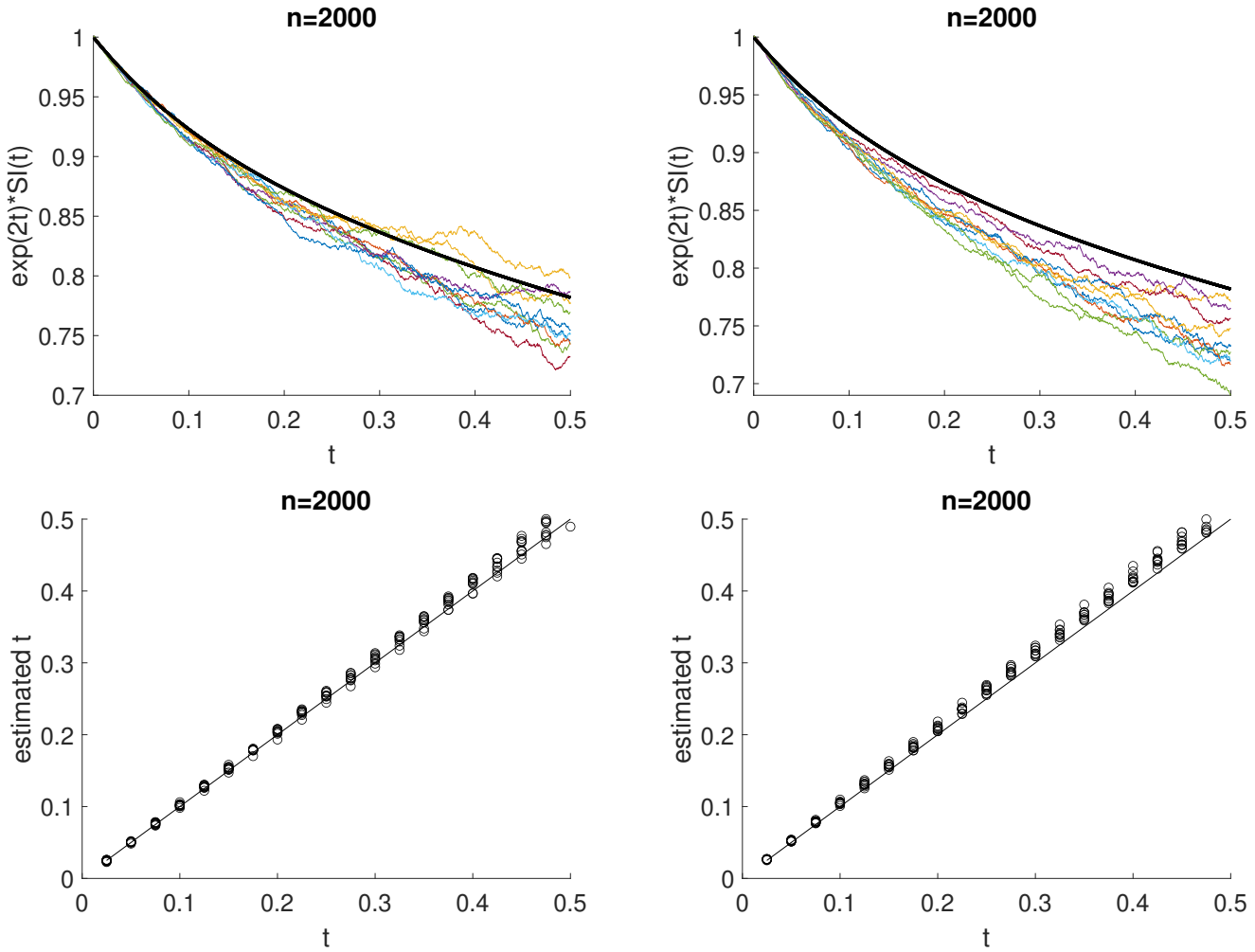


Fig. 3. **Simulation Results, Jump plus gain/loss:** Genome size $n = 2000$ Top: Estimated vs effective \hat{q}_k for gain/loss $p = 0.1, 0.2$. Bottom: Estimated vs effective \hat{t} for gain/loss $p = 0.1, 0.2$. Here the computation of the SI is performed using approach I, as described in the text.

364 base EggNOG (Huerta-Cepas et al., 2018) with 4.4M clusters of orthologous groups
 365 (COGs) (Tatusov et al., 2001). For each COG, EggNOG provides a flat ‘members’ file
 366 indicating the organisms that harbour this gene, along with its location in the genome.
 367 This allowed us to sort the genes by location along the genome. Within this representation,
 368 a genome is simply a list of COGs sorted by genome location, where the COG names are
 369 universal across all organisms. Hence, we can infer neighbourhood similarities across
 370 genomes and therefore the pairwise SI values between any two genomes which we then
 371 store in an $n \times n$ SI matrix. We set $k = 10$ which was found to be informative for these
 372 data (Shifman et al., 2013; Sevillya et al., 2019) and computed SI for all pairs of taxa. The
 373 crude SI values are strongly concentrated around 0.02, as shown in Fig. 6(R). In order to
 374 convert the SI values to a dissimilarity measure, we set $d_{SI} = 1 - SI$. Once a (pairwise)
 375 dissimilarity D matrix has been computed, we can then apply a distance-based
 376 phylogenetic method to estimate a tree T in which the leaves are labelled by the organisms
 377 under study.

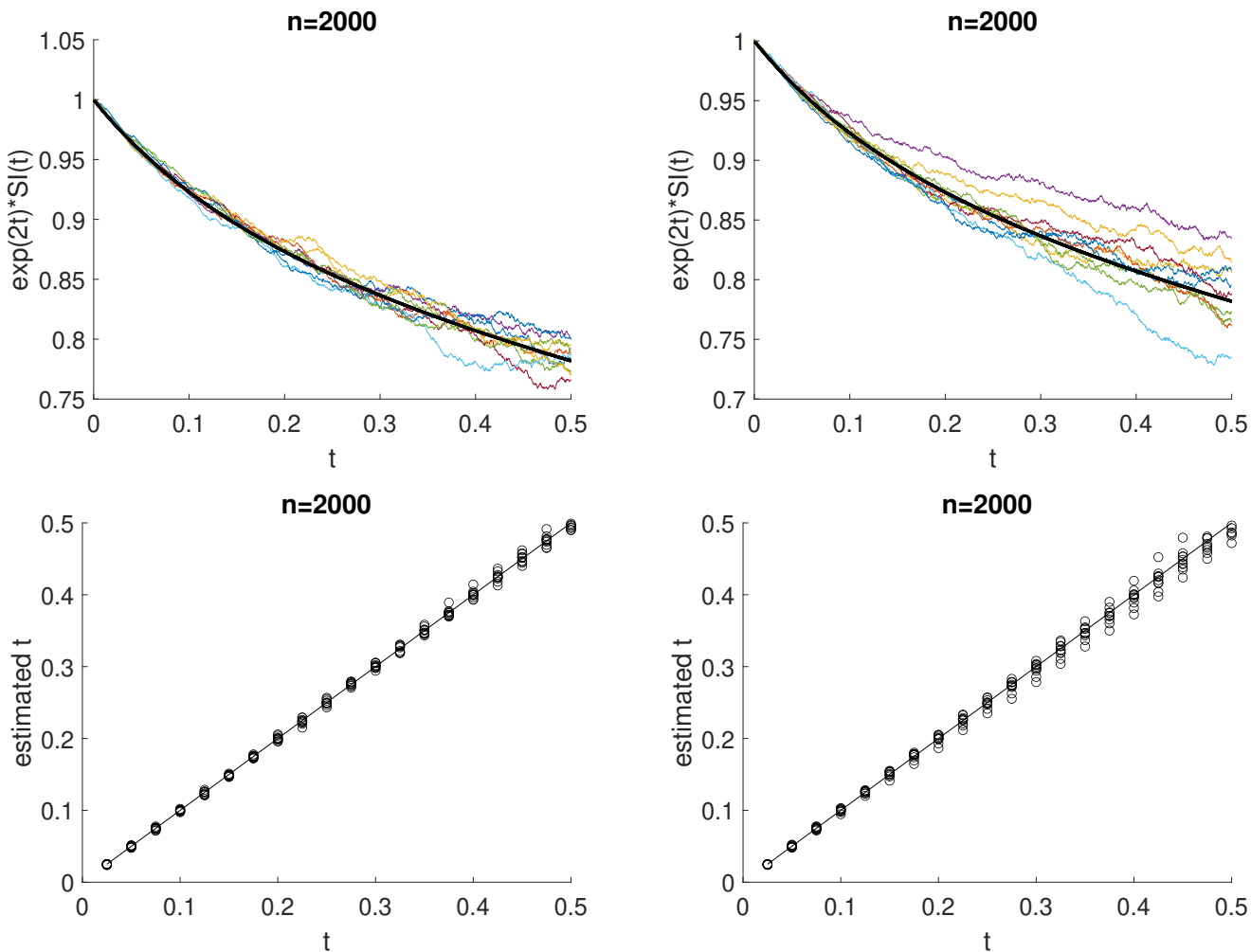


Fig. 4. **Simulation Results, Jump plus gain/loss:**. Genome size $n = 2000$ Top: Estimated vs effective \hat{q}_k for gain/loss $p = 0.1, 0.2$. Bottom: Estimated vs effective \hat{t} for gain/loss $p = 0.1, 0.2$. Here the computation of the SI is performed using approach II, as described in the text.

378 Path distances between the leaves of T , should approximate the distances in D . The
 379 most accurate algorithm for this task is the neighbour joining (NJ) algorithm (Saitou and
 380 Nei, 1987). Therefore, we used the program Neighbour from Phylip (Felsenstein, 1993) to
 381 construct a tree that we call the $1 - SI$ tree. Recall now that Eq. (8) was devised to
 382 “correct” the crude d_{SI} and provide a (provably) more reliable distance. Hence, we
 383 “corrected” the SI matrix accounting to Eq. (8) (specifically, finding \hat{t} by solving Eq. (8)
 384 for the appropriate SI value in the matrix) and then applied Neighbour to this matrix,
 385 yielding what we denote the *exact tree*. Finally, as in (Sevillya et al., 2019), we did not have
 386 an explicit expression for distance and were forced to develop a simulation-based heuristic,
 387 we also constructed the *heuristic tree* by using Formula (9) from (Sevillya et al., 2019).

388 EggNOG labels its organisms with the same taxon ID used by the NCBI taxonomy
 389 database (Federhen, 2011). This database is also furnished with taxonomic ranks in a
 390 child-to-parent relationship that we can use for our task. We therefore constructed a tree
 391 from this child–parent relationship. This NCBI tree spans about 1.2M organisms with

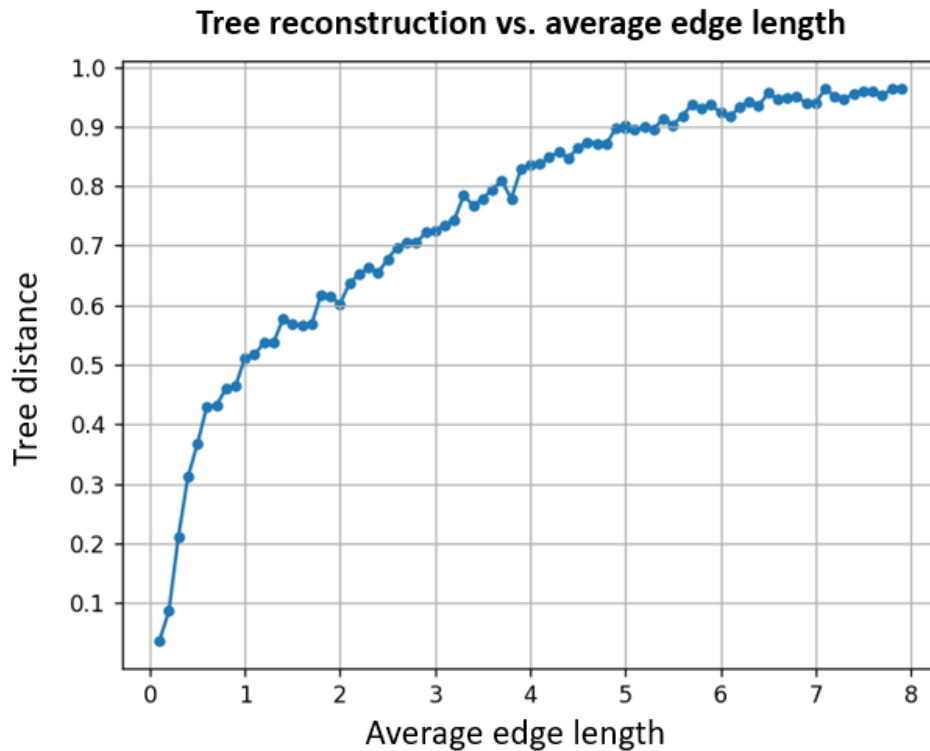


Fig. 5. **Simulation Results under a Tree Topology:** Genomes of size $n = 2000$ are evolved down a tree over 26 taxa under the same regime as in the above trials. The x and y axes are average edge length and normalized tree distance between the model, original tree, and the reconstructed tree.

392 maximum depth (i.e. ranks) of 39. We extracted the tree induced by EggNOG's 4445 taxa
 393 (this is done by removing the leaves that were not in the selection, and the paths leading
 394 solely to them) and used this tree as a reference tree, dubbed *the NCBI tree* (or
 395 *taxonomy*). The four trees appear in Fig. 7 in two formats - rectangular (L) and polar(R).
 396 As can be seen, the 1 – *SI* and the heuristic trees exhibit serious flaws we will elaborate on
 397 later. We wanted to measure the distance from each of the three SI- (or GD-, genome
 398 dynamic) based reconstructed trees to the reference NCBI tree. Again we used the
 399 Robinson–Foulds (RF) symmetric difference (Robinson and Foulds, 1981) tree metric. In
 400 presence of a reference, or a model, tree, the RF-distance can also be used to derive the
 401 false positive and false negative (FP, FN) rates. The relevant distances are presented in
 402 Figure 6(L). As can be seen, the exact tree from Eq. (8) is the most similar to the NCBI
 403 tree and the heuristic tree is the least similar.

404 The RF distance is very sensitive and uninformative for large trees (Siu-Ting et al.,
 405 2014). Hence, we adopted a coarser compatibility measure to allow a more intuitive
 406 assessment that can also detect differences between the two approaches, the sequence-based
 407 NCBI tree, and the other GD-based trees. We divided the NCBI reference tree into disjoint
 408 subtrees with sizes of between 80 and 800 taxa, resulting in 14 subtrees in total. This tree
 409 partitioning served as a *reference colouring* where each such NCBI subtree is mapped to a
 410 colour and all taxa (leaves) in this subtree attain that same colour (see the NCBI trees at
 411 the upper row of Fig. 7). As the NCBI Taxonomy provides classification to these families
 412 and genera represented by internal nodes in the NCBI tree, we could associate each such

413 subtree (or *clade*) with its corresponding evolutionary class. The classification is presented
 414 in Table 1.

Subtrees taxonomy

Color	Root ID	subtree description - NCBI Taxonomy
0	1117	Cyanobacteria (blue-green algae), phylum, cyanobacteria
1	72274	Pseudomonadales order, g-proteobacteria
2	91347	Enterobacteriales order, g-proteobacteria
3	135614	Xanthomonadales order, g-proteobacteria
4	135622	Alteromonadales order, g-proteobacteria
5	28211	Alphaproteobacteria class, a-proteobacteria
6	28216	Betaproteobacteria class, b-proteobacteria
7	68525	delta/epsilon subdivisions subphylum, proteobacteria
8	91061	Bacilli class, firmicutes
9	186801	Clostridia class, firmicutes
10	909932	Negativicutes class, firmicutes
11	68336	Bacteroidetes/Chlorobi group clade, bacteria
12	2037	Actinomycetales order, high G+C Gram-positive bacteria
13	544448	Tenericutes phylum, bacteria

Table 1. Subtrees taxonomy as provided by NCBI Taxonomy

415 The reference colouring of the NCBI tree allows us to measure this coloring in the
 416 other three SI-based trees as we describe next. In particular, as shown below, it allows us
 417 to detect incongruences between the SI-based and the sequence-based trees, that may
 418 suggest either misclassification or significant evolutionary events. Recall that the leaves in
 419 all trees are colored with the original color from the NCBI tree. Now, for a colour c , the
 420 c -subtree is defined as the minimal connected graph (subtree, in our case) containing all
 421 c -coloured leaves. Given a coloured tree (i.e. with some of the nodes coloured), such a
 422 colouring is said to be *convex* on that tree, if for every two colours c and c' , the c - and
 423 c' -coloured subtrees are disjoint (Moran and Snir, 2008, 2007). It is clear that the NCBI
 424 tree is convex, since the colouring is defined by this tree, i.e. for disjoint subtrees.
 425 Nevertheless, we aimed to test how far from convexity the NCBI colouring on the other
 426 trees is. There are rigorous definitions for the latter (the *recoloring distance* (Moran and
 427 Snir, 2008)); however, we used this approach to provide an intuitive and visual measure of
 428 compatibility, as demonstrated in Fig. 7.

429 As can be seen from the figure, all three trees maintained decent convexity under
 430 the NCBI colouring; however, it seems the exact tree has fewer violations than the
 431 heuristic and the $1 - SI$ trees. Fig. 7 also reveals major flaws in the heuristic and the
 432 $1 - SI$ approaches that are corrected by the exact approach. The $1 - SI$ approach takes
 433 crude values as the distances. These values are excessively concentrated around a tiny
 434 value of 0.02, causing severely distorted branch lengths, resulting in an artificially
 435 ultrametric tree with extremely short internal branches (third row in Fig. 7), which may
 436 disappear under bootstrapping, yielding a poorly resolved tree. Alternatively, the heuristic
 437 approach of (Sevillya et al., 2019), apart from achieving an outstandingly high RF
 438 distance, produces few exceptionally long branches non-proportional to the rest of the
 439 branches (left tree in the fourth row in Fig. 7). Hence, our real-data experiments showed
 440 that the theoretical conversion achieves its goal by producing a realistic distance, thereby

tree	Crude Robinson-Foulds				number of edges	common edges	% false positive	%false negative
NCBI Taxonomy	0	5516	4396	4358	1261			
Heuristic	5516	0	8884	8884	4443	94	98	92
1 – SI	4396	8884	0	2840	4443	654	85	48
Exact	4358	8884	2840	0	4443	673	85	47

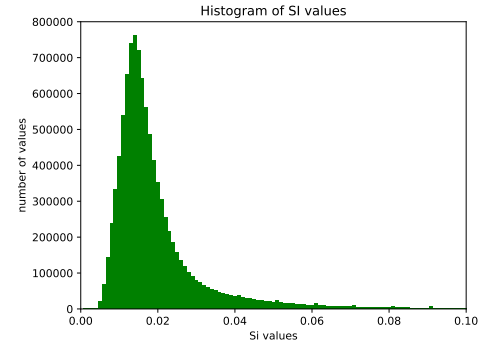


Fig. 6. **Left:** Robinson–Foulds distances. Pairwise RF distances between the four trees are depicted in columns 1–5. Column 6 contains the number of internal edges in each tree. Column 7 depicts the number of common edges with the NCBI tree, and Columns 8–9 depicts the rate of false positive and negative respectively (considering the NCBI tree as a model tree) **Right:** The distribution of pairwise SI values between all pairs from the 4445 EggNOG prokaryotic genomes.

441 correcting the severe flaws caused by the two more simplistic SI-based approaches.

442 As mentioned above, the incongruence in coloring between the two types of trees,
 443 GD-based (SI trees) and sequence-based (NCBI tree), may suggest closer scrutiny aiming
 444 at detecting genuine evolutionary phenomena. Thus, this procedure was followed. For each
 445 of the three SI-based trees and each of the 14 colors (corresponding to clades in the NCBI
 446 tree), we found the maximal subtree such that more than 90% of its leaves are colored in
 447 the desired color. For each such subtree, we counted the leaves colored in that color and
 448 their percentage, first of the total subtree size (number of leaves in the subtree), and next
 449 of the total number of leaves colored with that color. The exact numbers for each subtree
 450 appear in the table in supplementary text. While for some colors, all three SI-based tree
 451 are in nearly perfect agreement, conferring strong support in the NCBI classification, for
 452 other colors, all SI-based trees agree that the NCBI classification is incorrect. For example,
 453 the clades *Pseudomonadales*, *Alteromonadales*, the *delta/epsilon* subdivisions, and the
 454 *Bacilli* class, exhibit strong incongruence with the NCBI tree. In the supplementary text
 455 we provide details from the literature supporting the misclassification of these clades.
 456 Equivalently, this coloring provides hints to mislocation of specific taxa, a phenomenon
 457 referred to as *rogue taxa* (Smith, 2021). For example, all three SI-based trees allude to
 458 misclassification of taxa *gamma proteobacterium WG36*, *Gallaecimonas xiamenensis 3-C-1*,
 459 *Arsukibacterium perlucidum DSM 18276*, *Rheinheimera baltica DSM 14885*, and
 460 *Rheinheimera sp. A13L*, possibly mislocated based on their NCBI classification. Fuller
 461 details can be found in the supplementary text.

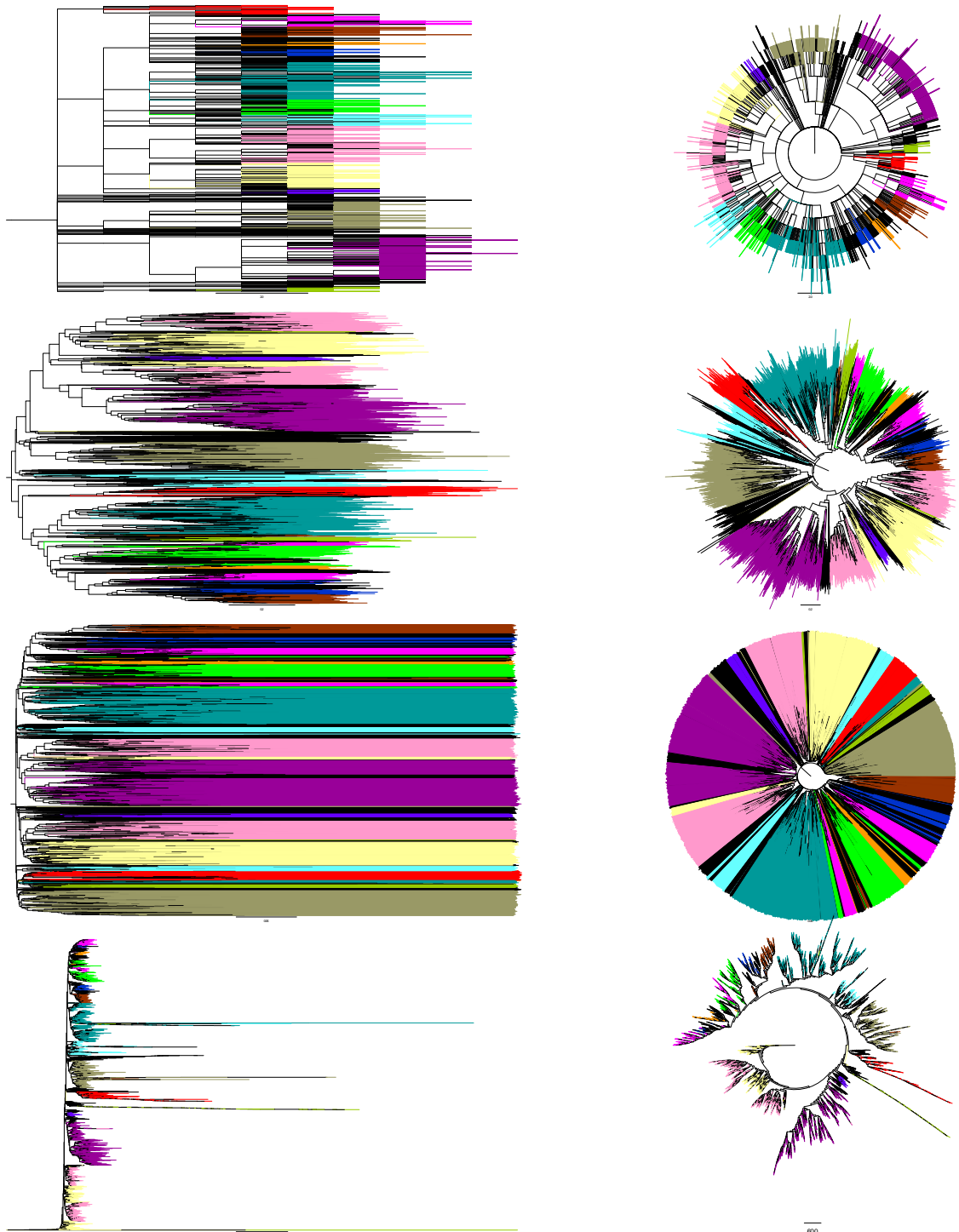


Fig. 7. **Coloured Trees:** Left: rectangular shape; right: polar shape. From the top: (1) The NCBI Taxonomy, (2) The exact SI tree, (3) the $1 - SI$ tree, (4) the heuristic exp. decay tree (the polar shape on right has log distances to accommodate the extremely long branches).

DISCUSSION

In this paper, we explored the consequences of modeling genome organisation as a continuous-time Markov process. Although the initial modelling was suggested recently, fundamental problems were left open, making it impossible to formally answer basic questions such as the time since divergence on a tree or the additivity of the synteny index as a phylogenetic marker. Here, we have advanced this front by applying mathematical tools from analysis and algebra to arrive at a rational function describing the transition probabilities, and the use of spectral theory and orthogonal polynomials, to prove the measure's consistency.

In the experimental realm, we demonstrated the accurate results provided by the new analytic expressions for real-life genome sizes and event rates. We have also extended the analytic model to account for realistic phenomena other than the Jump, and also on a tree structure. For the real data analysis, we built an ordered database of orthologous groups across 4445 prokaryotes, to which we applied our measure. To the best of our knowledge, there is no such database of this size in terms of orthologous groups or the number of taxa. Such a database could have multiple uses, apart from phylogenetics.

Applying our new measure to this database produced a tree that was in high accordance with the NCBI taxonomy for these organisms. Importantly, the new measure reconstructed realistic distances, as opposed to the previous measures, even the heuristic measure that was developed based on simulations. Reconstructing accurate distances has prime importance for establishing the Jump model as an underlying model of genome dynamics. Our results demonstrate that developing a distance measure, complementary to existing ones, is important for the sake of validating existing knowledge.

We expect that both the technique developed here for the modelling and the data resources will be instrumental in further analyses of other genome architectures such as operon and pseudogene formation.

While the Jump model is far from a precise description of the likely actual genome dynamics, its simplicity provides for analytical tractability. Future extensions of the model will account for more realistic scenarios including inversions of blocks of genes, duplications, and other events

DISCLOSURE STATEMENT

The authors have no conflicts of interest to declare. All co-authors have seen and agree with the contents of the manuscript and there is no financial interest to report. We certify that the submission is original work and is not under review at any other publication.

REFERENCES

- Adato, O., N. Ninyo, U. Gophna, and S. Snir. 2015. Detecting horizontal gene transfer between closely related taxa. *PLoS computational biology* 11:e1004408.
- Anderson, W. J. 2012. *Continuous-time Markov chains: An applications-oriented approach*. Springer Science & Business Media.
- Bansal, M. S., M. Kellis, M. Kordi, and S. Kundu. 2018. Ranger-dtl 2.0: rigorous reconstruction of gene-family evolution by duplication, transfer and loss. *Bioinformatics* 34:3214–3216.
- Baum, B. 1992. Combining trees as a way of combining data sets for phylogenetic inference. *Taxon* 41:3–10.

- 507 Biller, P., L. Guéguen, and E. Tannier. 2015. Moments of genome evolution by double
508 cut-and-join. *BMC bioinformatics* 16:S7.
- 509 Bininda-Emonds, O. 2004. Phylogenetic supertrees: Combining information to reveal the
510 Tree of Life. Springer.
- 511 Buneman, P. 1971. The recovery of trees from measures of dissimilarity. Pages 387–395 *in*
512 *Mathematics in the Archaeological and Historical Sciences* (F. Hodson, D. Kendall, and
513 P. Tautu, eds.). Edinburgh University Press, Edinburgh.
- 514 Che, D., G. Li, F. Mao, H. Wu, and Y. Xu. 2006. Detecting uber-operons in prokaryotic
515 genomes. *Nucleic acids research* 34:2418–2427.
- 516 Ciccarelli, F. D., T. Doerks, C. Von Mering, C. J. Creevey, B. Snel, and P. Bork. 2006.
517 Toward automatic reconstruction of a highly resolved tree of life. *science* 311:1283–1287.
- 518 Dalevi, D. and N. Eriksen. 2008. Expected gene-order distances and model selection in
519 bacteria. *Bioinformatics* 24:1332–1338.
- 520 Doolittle, W. F. 1999. Phylogenetic classification and the universal tree. *Science*
521 284:2124–2128.
- 522 Doyon, J.-P., C. Scornavacca, K. Y. Gorbunov, G. J. Szöllósi, V. Ranwez, and V. Berry.
523 2010. An efficient algorithm for gene/species trees parsimonious reconciliation with
524 losses, duplications and transfers. Pages 93–108 *in* RECOMB International Workshop on
525 Comparative Genomics Springer.
- 526 Federhen, S. 2011. The NCBI Taxonomy database. *Nucleic Acids Research* 40:D136–D143.
- 527 Felsenstein, J. 1978. Cases in which parsimony or compatibility methods will be positively
528 misleading. *Systematic zoology* 27:401–410.
- 529 Felsenstein, J. 1981. Evolutionary trees from dna sequences: a maximum likelihood
530 approach. *Journal of molecular evolution* 17:368–376.
- 531 Felsenstein, J. 1993. Phylip (phylogeny inference package), version 3.5 c .
- 532 Fitz Gibbon, S. T. and C. H. House. 1999. Whole genome-based phylogenetic analysis of
533 free-living microorganisms. *Nucleic acids research* 27:4218–4222.
- 534 Hannenhalli, S. and P. A. Pevzner. 1999. Transforming cabbage into turnip: polynomial
535 algorithm for sorting signed permutations by reversals. Pages 1–27 vol. 46 ACM.
- 536 Hasegawa, M., H. Kishino, and N. Saitou. 1991. On the maximum likelihood method in
537 molecular phylogenetics. *J Mol Evol* 32:443–445.
- 538 Huerta-Cepas, J., D. Szklarczyk, D. Heller, A. Hernandez-Plaza, S. K. Forslund, H. Cook,
539 D. R. Mende, I. Letunic, T. Rattei, L. Jensen, C. vonMering, and P. Bork. 2018. eggNOG
540 5.0: a hierarchical, functionally and phylogenetically annotated orthology resource based
541 on 5090 organisms and 2502 viruses. *Nucleic Acids Research* 47:D309–D314.
- 542 Jaccard, P. 1901. Étude comparative de la distribution florale dans une portion des alpes
543 et des jura. *Bull Soc Vaudoise Sci Nat* 37:547–579.
- 544 Koonin, E. V., K. S. Makarova, and L. Aravind. 2001. Horizontal gene transfer in
545 prokaryotes: quantification and classification. *Annual Reviews in Microbiology*
546 55:709–742.

- 547 Koonin, E. V., K. S. Makarova, and Y. I. Wolf. 2021. Evolution of microbial genomics:
548 Conceptual shifts over a quarter century. *Trends Microbiol* 29:582–592.
- 549 Lin, Y., F. Hu, J. Tang, and B. M. E. Moret. 2013. Maximum likelihood phylogenetic
550 reconstruction from high-resolution whole-genome data and a tree of 68 eukaryotes. *Pac*
551 *Symp Biocomput* Pages 285–296.
- 552 Liu, Y., P. Harrison, V. Kounin, and M. Gerstein. 2004. Comprehensive analysis of
553 pseudogenes in prokaryotes: widespread gene decay and failure of putative horizontally
554 transferred genes. *Genome Biol.* 5:R64.
- 555 Martinez-Gutierrez, C. A. and F. O. Aylward. 2021. Phylogenetic signal, congruence, and
556 uncertainty across bacteria and archaea. *Mol Biol Evol* 38:5514–5527.
- 557 Moran, S. and S. Snir. 2007. Efficient approximation of convex recolorings. *Journal of*
558 *Computer and System Sciences (JCSS)* 73:1078–1089 earlier version appeared in
559 *APPROX/RANDOM* 2005.
- 560 Moran, S. and S. Snir. 2008. Convex recolorings of strings and trees: Definitions, hardness
561 results and algorithms. *J. Comput. Syst. Sci.* 74:850–869.
- 562 Morel, B., A. M. Kozlov, A. Stamatakis, and G. J. Szöllösi. 2020. Generax: A tool for
563 species-tree-aware maximum likelihood-based gene family tree inference under gene
564 duplication, transfer, and loss. *Mol Biol Evol* 37:2763–2774.
- 565 Morel, B., P. Schade, S. Lutteropp, T. A. Williams, G. J. Szllsi, and A. Stamatakis. 2022.
566 SpeciesRax: A Tool for Maximum Likelihood Species Tree Inference from Gene Family
567 Trees under Duplication, Transfer, and Loss. *Molecular Biology and Evolution* 39
568 msab365.
- 569 Nakhleh, L., D. Ruths, and L.-S. Wang. 2005. Riata-hgt: a fast and accurate heuristic for
570 reconstructing horizontal gene transfer. Pages 84–93 *in* *International Computing and*
571 *Combinatorics Conference*.
- 572 Ochman, H., J. G. Lawrence, and E. A. Groisman. 2000. Lateral gene transfer and the
573 nature of bacterial innovation. *nature* 405:299.
- 574 Pang, T. Y. and M. J. Lercher. 2019. Each of 3,323 metabolic innovations in the evolution
575 of *E. coli* arose through the horizontal transfer of a single DNA segment. *Proceedings of*
576 *the National Academy of Sciences of the United States of America* 116:187–192.
- 577 Puigbò, P., A. E. Lobkovsky, D. M. Kristensen, Y. I. Wolf, and E. V. Koonin. 2014.
578 Genomes in turmoil: quantification of genome dynamics in prokaryote supergenomes.
579 *BMC biology* 12:66.
- 580 Pybus, O. G. and A. Rambaut. 2009. Evolutionary analysis of the dynamics of viral
581 infectious disease. *Nature Reviews Genetics* 10:540–550.
- 582 Ragan, M. 1992. Matrix representation in reconstructing phylogenetic-relationships among
583 the eukaryotes. *Biosystems* 28:47–55.
- 584 Rajendhran, J. and P. Gunasekaran. 2011. Microbial phylogeny and diversity: small
585 subunit ribosomal rna sequence analysis and beyond. *Microbiol Res* 166:99–110.
- 586 Robinson, D. F. and L. R. Foulds. 1981. Comparison of phylogenetic trees. *Mathematical*
587 *biosciences* 53:131–147.

- 588 Rogozin, I., K. Makarova, J. Murvai, E. Czabarka, Y. Wolf, R. Tatusov, L. Szekely, and
589 E. Koonin. 2002. Connected gene neighborhoods in prokaryotic genomes. *Nucleic Acids*
590 *Res* 30:2212–2223.
- 591 Saitou, N. and M. Nei. 1987. The neighbor-joining method: a new method for
592 reconstructing phylogenetic trees. *Molecular Biology and Evolution* 4:406–425.
- 593 Sankoff, D. 1992. Edit distance for genome comparison based on non-local operations.
594 Pages 121–135 *in* Annual Symposium on Combinatorial Pattern Matching Springer.
- 595 Sankoff, D. and J. H. Nadeau. 1996. Conserved synteny as a measure of genomic distance.
596 *Discrete applied mathematics* 71:247–257.
- 597 Schoch, C. L., S. Ciuffo, M. Domrachev, C. L. Hottel, S. Kannan, R. Khovanskaya,
598 D. Leipe, R. Mcveigh, K. O'Neill, B. Robbertse, S. Sharma, V. Soussov, J. P. Sullivan,
599 L. Sun, S. Turner, and I. Karsch-Mizrachi. 2020. NCBI Taxonomy: a comprehensive
600 update on curation, resources and tools. *Database* 2020 baaa062.
- 601 Schönknecht, G., A. P. M. Weber, and M. J. Lercher. 2014. Horizontal gene acquisitions by
602 eukaryotes as drivers of adaptive evolution. *BioEssays* 36:9–20 ISBN: 1521-1878
603 (Electronic)\r0265-9247 (Linking).
- 604 Semple, C. and M. Steel. 2003. *Phylogenetics* vol. 24. Oxford University Press on Demand.
- 605 Serdoz, S., A. Egri-Nagy, J. Sumner, B. R. Holland, P. D. Jarvis, M. M. Tanaka, and A. R.
606 Francis. 2017. Maximum likelihood estimates of pairwise rearrangement distances.
607 *Journal of theoretical biology* 423:31–40.
- 608 Sevillya, G., D. Doerr, Y. Lerner, J. Stoye, M. Steel, and S. Snir. 2019. Horizontal Gene
609 Transfer Phylogenetics: A Random Walk Approach. *Molecular Biology and Evolution*
610 37:1470–1479.
- 611 Sevillya, G. and S. Snir. 2019. Synteny footprints provide clearer phylogenetic signal than
612 sequence data for prokaryotic classification. *Molecular phylogenetics and evolution*
613 136:128–137.
- 614 Shifman, A., N. Ninyo, U. Gophna, and S. Snir. 2013. Phylo si: a new genome-wide
615 approach for prokaryotic phylogeny. *Nucleic acids research* 42:2391–2404.
- 616 Siu-Ting, K., D. Pisani, C. J. Creevey, and M. Wilkinson. 2014. Concatabominations:
617 Identifying Unstable Taxa in Morphological Phylogenetics using a Heuristic Extension to
618 Safe Taxonomic Reduction. *Systematic Biology* 64:137–143.
- 619 Sjöstrand, J., A. Tofigh, V. Daubin, L. Arvestad, B. Sennblad, and J. Lagergren. 2014. A
620 bayesian method for analyzing lateral gene transfer. *Systematic biology* 63:409–420.
- 621 Smith, M. R. 2021. Using Information Theory to Detect Rogue Taxa and Improve
622 Consensus Trees. *Systematic Biology* 71:1088–1094.
- 623 Snel, B., P. Bork, and M. A. Huynen. 1999. Genome phylogeny based on gene content.
624 *Nature genetics* 21:108.
- 625 Stolzer, M., H. Lai, M. Xu, D. Sathaye, B. Vernet, and D. Durand. 2012. Inferring
626 duplications, losses, transfers and incomplete lineage sorting with nonbinary species
627 trees. *Bioinformatics* 28:i409–i415.
- 628 Strimmer, K. and V. Moulton. 2000. Likelihood analysis of phylogenetic networks using
629 directed graphical models. *Mol Biol Evol* 17:875–881.

- 630 Szöllősi, G. J., E. Tannier, N. Lartillot, and V. Daubin. 2013. Lateral gene transfer from
631 the dead. *Systematic biology* 62:386–397.
- 632 Tatusov, R. L., D. A. Natale, I. V. Garkavtsev, T. A. Tatusova, U. T. Shankavaram, B. S.
633 Rao, B. Kiryutin, M. Y. Galperin, N. D. Fedorova, and E. V. Koonin. 2001. The cog
634 database: new developments in phylogenetic classification of proteins from complete
635 genomes. *Nucleic acids research* 29:22–28.
- 636 Tekaiia, F. and B. Dujon. 1999. Pervasiveness of gene conservation and persistence of
637 duplicates in cellular genomes. *Journal of molecular evolution* 49:591–600.
- 638 Wang, L.-S. and T. Warnow. 2001. Estimating true evolutionary distances between
639 genomes. Pages 637–646 *in* Proceedings of the thirty-third annual ACM symposium on
640 Theory of computing ACM.
- 641 Yancopoulos, S., O. Attie, and R. Friedberg. 2005. Efficient sorting of genomic
642 permutations by translocation, inversion and block interchange. *Bioinformatics*
643 21:3340–3346.
- 644 Yang, Z. 1996. Maximum-likelihood models for combined analyses of multiple sequence
645 data. *Journal of Molecular Evolution* 42:587–596.
- 646 Zhao, T., A. Zwaenepoel, J.-Y. Xue, S.-M. Kao, Z. Li, M. E. Schranz, and Y. Van de Peer.
647 2021. Whole-genome microsynteny-based phylogeny of angiosperms. *Nat Commun*
648 12:3498.

APPENDIX: MATHEMATICAL PROOFS

Proof for Theorem 2

We now provide a proof for Theorem 2 provided in the main text. We first repeat the theorem.

Theorem 2:

$$p_{i,j}(t) = \frac{1}{(t+1)^{i+j-1}} \cdot \sum_{\ell=1}^{\min(i,j)} \frac{(i+j-\ell-1)!}{(i-\ell)!(j-\ell)!(\ell-1)!} (1-t^2)^{\ell-1} t^{i+j-2\ell}. \quad (9)$$

Proof: This result follows from some general results for birth–death processes (refer to (Anderson, 2012) for these). A simple change in notation will be needed, since the results of (Anderson, 2012) involve a birth–death process that is defined on the non-negative integers, whereas our process above is defined on the positive integers. We therefore define $Y_t = X_t - 1$, so that the process Y_t satisfies an instance of the general birth–death process described by:

$$\begin{aligned} \mathbb{P}(Y_{t+\delta} = i+1 | Y_t = i) &= \lambda_i \delta + o(\delta), \quad i \geq 0 \\ \mathbb{P}(Y_{t+\delta} = i-1 | Y_t = i) &= \mu_i \delta + o(\delta), \quad i \geq 1 \end{aligned}$$

where in our case:

$$\lambda_i = i+1, \mu_i = i. \quad (10)$$

At the heart of the spectral theory of birth–death processes is the Karlin-McGregor representation of the state transition probabilities ((Anderson, 2012), Ch. 8, Theorem 2.1):

$$\mathbb{P}(Y_t = j | Y_0 = i) = \pi_j \int_0^\infty e^{-tx} Q_i(x) Q_j(x) d\psi(x), \quad (11)$$

where $d\psi(x)$ is a measure on $[0, \infty)$, known as the spectral measure, $\{Q_i(x)\}_{i=0}^\infty$ is a sequence of polynomials, orthogonal with respect to the measure $d\psi$, and $\pi_j = \prod_{k=0}^{j-1} \frac{\lambda_k}{\mu_{k+1}}$.

In the particular case where the birth–death process is given by (10), we have ((Anderson, 2012), Ch.8, Eq. 4.14):

$$d\psi(x) = e^{-x} dx, \quad (12)$$

$$\pi_j = 1, \quad j \geq 0, \quad (13)$$

and the polynomials $Q_i(x)$ are the Laguerre polynomials defined by ((Anderson, 2012), Ch.8, Eq. 4.12)

$$Q_i(x) = {}_1F_1(-m; 1, x) = \sum_{k=0}^i \frac{(-i)_k}{k!^2} \cdot x^k$$

where ${}_1F_1$ is the confluent hypergeometric function, and $(-i)_k = (-i)(-i+1)\cdots(-i+k-1)$. We also have the relation ((Anderson, 2012), Ch.8, eq. 4.15)

$$\int_0^\infty e^{-sx} Q_i(x) Q_j(x) dx = \frac{(i+j)!}{i!j!} \cdot \frac{(s-1)^{i+j}}{s^{i+j+1}} \cdot {}_2F_1\left(-i, -j; -i-j; \frac{s(s-2)}{(s-1)^2}\right), \quad (14)$$

where ${}_2F_1$ is the Gaussian hypergeometric function, defined by:

$${}_2F_1\left(-i, -j; -i-j; \frac{s(s-2)}{(s-1)^2}\right) = \sum_{k=0}^\infty \frac{(-i)_k (-j)_k}{(-i-j)_k} \cdot \left(\frac{s(s-2)}{(s-1)^2}\right)^k. \quad (15)$$

Using (12),(13),(14),(15) with $s = t + 1$, (11) leads to ((Anderson, 2012), Ch. 8, eq. 4.28):

$$\begin{aligned}
\mathbb{P}(Y_t = j | Y_0 = i) &= \pi_j \int_0^\infty e^{-tx} Q_i(x) Q_j(x) e^{-x} dx = \int_0^\infty e^{-x(t+1)} Q_i(x) Q_j(x) dx \\
&= \frac{(i+j)!}{i!j!} \cdot \frac{t^{i+j}}{(t+1)^{i+j+1}} \cdot {}_2F_1 \left(-i, -j; -i-j; \frac{t^2-1}{t^2} \right) \\
&= \frac{(i+j)!}{i!j!} \cdot \frac{t^{i+j}}{(t+1)^{i+j+1}} \sum_{k=0}^\infty \frac{(-i)_k (-j)_k}{(-i-j)_k k!} \cdot \left(\frac{t^2-1}{t^2} \right)^k \\
&= \frac{(i+j)!}{i!j!} \cdot \frac{t^{i+j}}{(t+1)^{i+j+1}} \sum_{k=0}^{\min(i,j)} \frac{i!j!(i+j-k)!(-1)^k}{(i-k)!(j-k)!(i+j)!k!} \cdot \left(\frac{t^2-1}{t^2} \right)^k \\
&= \frac{t^{i+j}}{(t+1)^{i+j+1}} \sum_{k=0}^{\min(i,j)} \frac{(i+j-k)!}{(i-k)!(j-k)!k!} \cdot \left(\frac{1-t^2}{t^2} \right)^k.
\end{aligned}$$

Therefore, going back from the process Y_t to the process X_t , we have

$$\begin{aligned}
p_{i,j}(t) &= p_{ij}(t) = \mathbb{P}(X_t = j | X_0 = i) = p_{ij}(t) = \mathbb{P}(Y_t = j+1 | Y_0 = i+1) \\
&= \frac{t^{i+j-2}}{(t+1)^{i+j-1}} \sum_{k=0}^{\min(i,j)-1} \frac{(i+j-k-2)!}{(i-1-k)!(j-1-k)!k!} \cdot \left(\frac{1-t^2}{t^2} \right)^k \\
&= \frac{1}{(t+1)^{i+j-1}} \cdot \sum_{\ell=1}^{\min(i,j)} \frac{(i+j-\ell-1)!}{(i-\ell)!(j-\ell)!(\ell-1)!} (1-t^2)^{\ell-1} t^{i+j-2\ell}.
\end{aligned}$$

666

□

Proof for Theorem 3

667

Theorem 3:

668

$$q_k(t) = \frac{1}{k} \sum_{\ell=0}^{k-1} \sum_{i=0}^{k-\ell-1} \sum_{j=0}^{k-\ell-1} \frac{(i+j+\ell)!}{i!j!\ell!} t^{i+j} (t+1)^{-i-j-2\ell-1} (1-t^2)^\ell. \quad (16)$$

Proof. Summing the expressions for $p_{ij}(t)$ we get

$$\begin{aligned}
q_k(t) &= \frac{1}{k} \sum_{i=1}^k \sum_{j=1}^k p_{ij}(t) = \sum_{i=1}^k \sum_{j=1}^k \frac{1}{(t+1)^{i+j-1}} \cdot \sum_{\ell=1}^{\min(i,j)} \frac{(i+j-\ell-1)!}{(i-\ell)!(j-\ell)!(\ell-1)!} (1-t^2)^{\ell-1} t^{i+j-2\ell} \\
&= \frac{1}{k} (t+1) \sum_{\ell=1}^k \frac{1}{(\ell-1)!} t^{-2\ell} (1-t^2)^{\ell-1} \sum_{i=\ell}^k \frac{t^i}{(t+1)^i} \frac{1}{(i-\ell)!} \sum_{j=\ell}^k \frac{t^j}{(t+1)^j} \frac{(i+j-\ell-1)!}{(j-\ell)!} \\
&= \frac{1}{k} (t+1) \sum_{\ell=1}^k \frac{1}{(\ell-1)!} t^{-2\ell} (1-t^2)^{\ell-1} \sum_{i=\ell}^k \frac{t^i}{(t+1)^i} \frac{1}{(i-\ell)!} \sum_{j=0}^{k-\ell} \frac{t^{j+\ell}}{(t+1)^{j+\ell}} \frac{(i+j-1)!}{j!}
\end{aligned}$$

$$\begin{aligned}
&= \frac{1}{k} \frac{1}{(t+1)^{2k-1}} \sum_{\ell=0}^{k-1} \frac{1}{\ell!} (1-t^2)^\ell \sum_{i=0}^{k-\ell-1} \sum_{j=0}^{k-\ell-1} \frac{(i+j+\ell)!}{i!j!} t^{i+j} (t+1)^{2k-i-j-2\ell-2} \\
&= \frac{1}{k} \sum_{\ell=0}^{k-1} \sum_{i=0}^{k-\ell-1} \sum_{j=0}^{k-\ell-1} \frac{(i+j+\ell)!}{i!j!\ell!} t^{i+j} (t+1)^{-i-j-2\ell-1} (1-t^2)^\ell.
\end{aligned}$$

669

□

Proof for Theorem 4

670

671 **Theorem 4:** The function $q_k(t)$ is monotone decreasing on $[0, \infty)$.
Proof. Using the representation given by Eq. (11) we have:

$$p_{ij}(t) = \int_0^\infty e^{-tx} Q_{i-1}(x) Q_{j-1}(x) d\psi(x).$$

This implies that

$$\begin{aligned}
q_k(t) &= \frac{1}{k} \sum_{i=1}^k \sum_{j=1}^k p_{ij}(t) = \frac{1}{k} \sum_{i=1}^k \sum_{j=1}^k \int_0^\infty e^{-tx} Q_{i-1}(x) Q_{j-1}(x) d\psi(x) \\
&= \frac{1}{k} \int_0^\infty e^{-tx} \sum_{i=1}^k \sum_{j=1}^k Q_{i-1}(x) Q_{j-1}(x) d\psi(x) \\
&= \frac{1}{k} \int_0^\infty e^{-tx} \left(\sum_{i=1}^k Q_{i-1}(x) \right)^2 d\psi(x)
\end{aligned}$$

Therefore, by differentiating the above with respect to t we obtain:

$$q'_k(t) = -\frac{1}{k} \int_0^\infty e^{-tx} x \left(\sum_{i=1}^k Q_{i-1}(x) \right)^2 d\psi(x) < 0,$$

672

since the integrand is positive. This establishes that $q_k(t)$ is monotone decreasing. □