

Introduction

Through rapid evolution, human influenza viruses are able to evade host immunity in populations around the globe. In addition to mutation, reassortment of the different segments of influenza viruses provides an important source of viral diversity (Steel and Lowen, 2014). If a cell is infected by more than one virus, progenitor viruses can carry segments from more than one parent (McDonald et al., 2016). with the exception of accidental release of antigenically lagged human influenza viruses (Nakajima, Desselberger, and Palese, 1978), reassortment remains the sole documented mechanism for generating pandemic influenza strains (e.g (Smith, Bahl, et al., 2009; Smith, Vijaykrishna, et al., 2009; Guan et al., 2010)).

To characterize such events, tanglegrams, comparison between tree heights (Westgeest et al., 2014; Dudas, Bedford, et al., 2014), or ancestral state reconstructions (Lu, Lycett, and Brown, 2014) are typically deployed. These approaches identify discordance between different segment tree topologies or differences in pairwise distances between isolates across segment trees. Tanglegrams in particular require a substantial amount of subjectivity and have been described as potentially misleading (De Vienne, 2018).

While the reassortment process has been intensively studied (e.g (Nelson et al., 2008; Westgeest et al., 2014; Dudas, Bedford, et al., 2014; Lu, Lycett, and Brown, 2014)), there is currently no explicit model based inference approach available. We address this void by introducing a coalescent-based model of the reassortment of viral lineages. In this phylogenetic network model, ancestral lineages carry genome segments, of which only a subset may be ancestral to sampled viral genomes. As in a normal coalescent process, network lineages coalesce (merge) with each other backwards in time at a rate inversely proportional to the effective population size. We model reassortment (splitting) events as a result of a constant-rate Poisson process on network lineages. At such a splitting event, the ancestry of segments on the original lineage diverges, with a random subset following each new lineage. We thus explicitly model reassortment networks and the embedding of segment trees within these, allowing us to infer these entities from available sequence data.

In order to perform inference under such a model, the reassortment network and the embedding of each segment tree within that network must be jointly inferred. This is similar to the well-known and challenging problem of inferring ancestral recombination graphs (ARGs), with the difference that segments in our model have fixed boundaries, but no defined ordering. While many approaches to inferring ARGs exist, some are restricted to tree-based networks (Didelot et al., 2010; Vaughan, Welch, et al., 2017), meaning that the networks consist of a base tree where recombination edges always attach to edges on the base tree. Other approaches (e.g (M. D. Rasmussen et al., 2014)) rely on approximations (McVean and Cardin, 2005) and are not applicable to the reassortment model due to the aforementioned lack of segment ordering. Completely general inference methods exist (Erik W. Bloomquist and Marc A. Suchard, 2010), but these are again not directly applicable to modelling reassortment and furthermore tend to be highly computationally demanding.

Here we introduce a novel Markov Chain Monte Carlo (MCMC) approach to jointly sample networks and the embedding of segment trees within those networks, without any approximations involved. This approach allows us to perform joint inference of the reassortment network,

87 the phylogenetic trees of each segment, the reassortment and coalescent
88 rates with the evolutionary models and its parameters.

89 We first show that this approach is able to retrieve reassortment rates,
90 effective population sizes and reassortment events from simulated data.
91 Secondly, we discuss how a lack of genetic information influences the in-
92 ference of these parameters. Thirdly, we show how using the coalescent
93 with reassortment can influence the inference of effective population sizes,
94 as well as evolutionary rates. We then apply this approach to study re-
95 assortment rates patterns across different influenza subtypes. Finally, we
96 study how reassortment rates differ on fit and unfit edges of these reas-
97 sortment networks.

98 **Inference of effective population sizes and reassort-** 99 **ment rates are reliably inferred from genetic se-** 100 **quence data**

101 In order to test our ability to infer effective population sizes and reas-
102 sortment rates from genetic sequences, we performed a well calibrated
103 simulation study. To do so, we first sampled random effective population
104 sizes from a log normal distribution (mean 5 and standard deviation 0.5)
105 and reassortment rates from another log normal distribution (mean 0.2
106 and standard deviation 0.5). We then sampled the sampling times of 100
107 taxa, each with 4 segments, from a uniform distribution between 0 and
108 20. We next simulated reassortment networks alongside the embedding of
109 the segment trees using these parameters. For each segment tree, we next
110 simulated genetic sequences by using the JC69 substitution model (Jukes
111 and Cantor, 1969) with an evolutionary rate of 5×10^{-3} per site and
112 year. Each segment thereby consisted of 1000 independently evolving nu-
113 cleotides. In order to study the effect of reducing the amount of genetic
114 information, we additionally considered the scenario where all segments
115 had an evolutionary rate of 5×10^{-4} per site and year. Using our MCMC
116 approach we then inferred the reassortment network, segment tree embed-
117 ding, effective population sizes and reassortment rates from these genetic
118 sequences.

119 The results shown in figure 1A,B, indicate that we are able to correctly
120 retrieve effective population sizes and reassortment rates from simulated
121 genetic sequences. Effective population sizes are estimated more precisely
122 than reassortment rates, which is expected considering that there are typ-
123 ically many more coalescent events in a network than reassortment events.
124 Lower evolutionary rates do not greatly decrease our ability to infer effec-
125 tive population sizes and reassortment rates (see figure S1).

126 To test how well true reassortment events are recovered, we com-
127 puted the probability of observing exactly the same reassortment events
128 as present in the true (simulated) network. We considered to reassort-
129 ment events to be the same if the sub-tree of each segment below that
130 node is the same and if the relative direction of each segment at the reas-
131 sortment event is exactly the same (see Methods, Network Summary).
132 This constitutes a stringent definition of two reassortment events being
133 the same.

134 As shown in figure 1C, reassortment events are well supported, particu-
135 larly with increasing reassortment distance. The reassortment distance
136 denotes how much independent evolution happened on the two parent
137 viruses of the reassortment event (see Methods, Reassortment Distance).

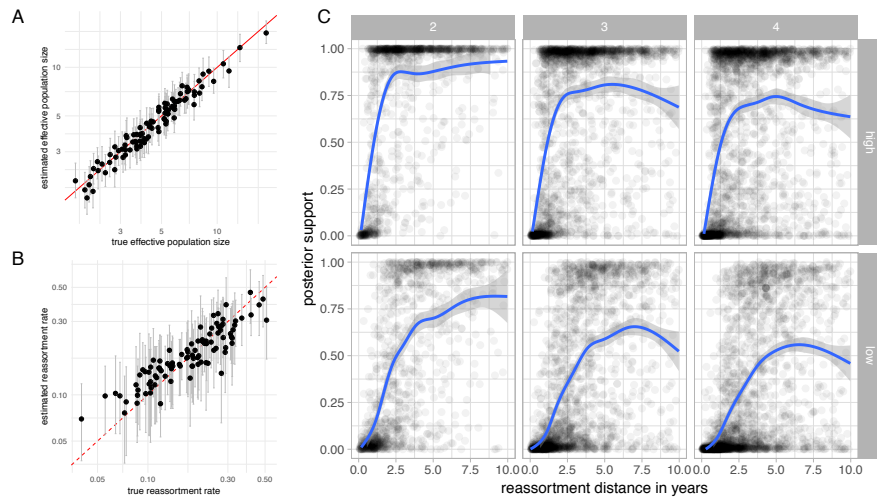


Figure 1: Estimates of effective population sizes and reassortment rates from simulated genetic sequences. **A** Estimated effective population sizes and 95% confidence intervals (y-axis) vs. simulated effective population sizes on the x-axis. **B** Estimated reassortment rates and 95% confidence (y-axis) vs. simulated reassortment rates on the x-axis. **C** Posterior support for true reassortment events (y-axis) given the reassortment distance (x-axis). Inference of reassortment networks from sequences simulated with a evolutionary rate of 5×10^{-3} mutations per site and year (top row) and 5×10^{-4} mutations per site and year (bottom row). From left to right, the reassortment events are for networks with 2,3 and 4 segments.

138 This is particularly true when we only look at reassortment events between
 139 pairs of segments and drops when we look at 3 or 4 segments. This
 140 decrease is driven by our definition that two reassortment events are only
 141 the same if all segments reassort in the same relative direction at the same
 142 time with exactly the same clade below the segment trees; a requirement
 143 that becomes harder to satisfy as the number of segments increases. As
 144 expected for methods that correctly take into account uncertainty, the
 145 posterior support decreases when lower evolutionary rates are used to
 146 simulate the sequences of the segments.

147 **Joint inference from full genomes increases precision in dating nodes**

149 We compared the internal node ages inferred using the coalescent with
 150 reassortment to ages inferred under the assumption that all segments
 151 evolved independently under the standard coalescent model. To do this,
 152 we first compiled datasets of several human influenza A subtypes, as well
 153 as influenza B (details in the Materials and Methods). From each of these
 154 we generated three datasets consisting of a random sample of sequences.

155 We then analysed each of these sub-sampled datasets once using the
 156 coalescent with reassortment and once using a normal coalescent prior
 157 with shared effective population size across all segments, but assuming
 158 that each segment evolved independently. We first computed the 95%

159 highest posterior density (HPD) interval of node heights for each clade
160 that was supported by both approaches with a posterior probability of
161 more than 0.5. We then normalized the difference between the lower and
162 upper bound of the 95% HPD interval, by the median node height estimate
163 to get the relative width of the HPD interval for each clade. As shown in
164 figure 2A, using the coalescent with reassortment reduces the uncertainty
165 of node height estimates of segment tree nodes by 34% for p2009 like
166 H1N1) up to 50% for influenza B.

167 Next we computed the distribution of clade supports for clades repre-
168 sented in the MCC trees inferred using the two approaches. As shown in
169 figure 2B, segment tree clades are much better resolved when using the
170 coalescent with reassortment for all datasets.

171 We then compared the effective population sizes and evolutionary inferred
172 using the two approaches. The coalescent with reassortment infers
173 higher effective population sizes for all datasets (see figure 2C). This also
174 influences the inferred clock rates, since lower effective population sizes
175 put stronger weight on shorter branches and therefore larger clock rates
176 (see figure 2C). We explain this discrepancy as follows. Coalescent events
177 closer to the tips are more likely between lineages that are for example
178 geographically more closely related and can be assumed to occur rapidly
179 and provide information about low effective population size values. Coales-
180 cent events deeper in the tree on the other hand are more representative
181 of those between geographically more separated lineages. These events
182 therefore provide information about larger effective population sizes. Co-
183 alescent events across different segments that occur close to the tips are
184 less likely to have encountered reassortment events. In the coalescent with
185 reassortment, they are therefore interpreted as one event, whereas in the
186 coalescent with independent segments, they are interpreted as eight. Co-
187 alescent events deeper in the tree are more likely between lineages that
188 encountered reassortment events and are therefore more likely to provide
189 independent information about the population process. The coalescent
190 with independent segments assumes that all coalescent events provide the
191 same amount of information about the population process and will conse-
192 quently favour information about the population process closer to the
193 tips. This leads to differences in the estimated effective population sizes
194 which then leads to differences in the estimated clock rates.

195 We also compared the performance of the two approaches by infer-
196 ring tip dates (Dudas and Bedford, 2019). The tips (leaf nodes) are the
197 only nodes in the trees or network for which we can actually presume to
198 know the true age, which is set by the sample collection time. To compare
199 the two approaches, we compiled 1000 smaller influenza A/H3N2 datasets
200 each composed of 20 genomes. Of those datasets, 500 were randomly sam-
201 pled from an interval of 2 years between 1995 and 2019. The remaining 500
202 datasets were assembled using a random sampling interval of 10 years be-
203 tween 1995 and 2019. From each of these datasets we randomly selected
204 a single genome and inferred its sampling time using both approaches,
205 conditional on the sampling times of the remaining genomes.

206 The 95% HPD of the sample time posteriors under the coalescent with
207 reassortment contains the true sampling time interval in 91% of cases
208 for the 2 year sampling interval and in 89% for the 10 year sampling
209 interval (see figure S2). On the other hand, the 95% HPD of the sample
210 time posteriors generated by the independent segment coalescent model
211 contains the true sampling time in only 68% (2 year sampling interval)
212 and 77% (10 year sampling interval) of cases.

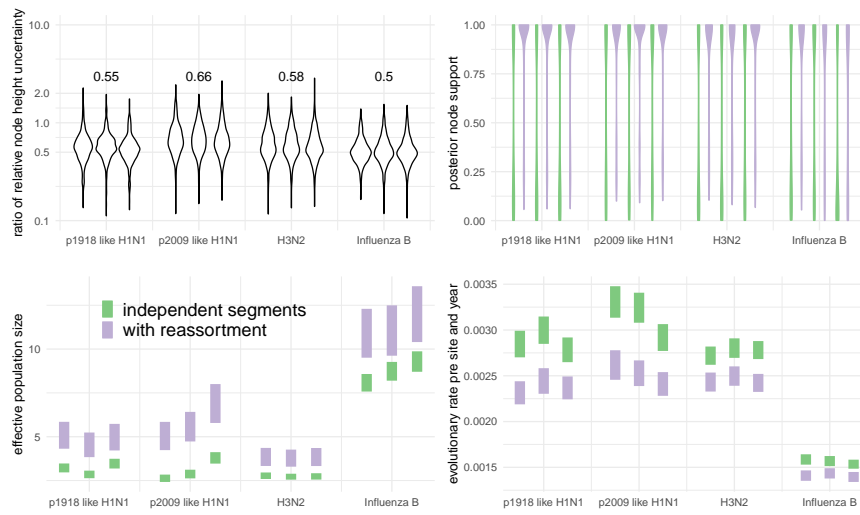


Figure 2: Comparison of estimates between the coalescent with reassortment and assuming that each segment codes for an independent realization of the same coalescent process. **A** Comparison of the relative width of the 95% HPD interval of segment tree node heights. The vertical axis shows the distribution of ratios of the relative width of the 95% HPD intervals of the coalescent with reassortment over the coalescent assuming independent segment evolution. The values show the median reduction in node height uncertainty when using the coalescent with reassortment over the coalescent with independent segments. **B** Comparison between the distribution of posterior clade support of segment trees found the maximum clade credibility segment trees. **C** Comparison between the inferred effective population sizes. When assuming each segment is an independent realization of the same coalescent process, the effective population sizes are inferred to be much smaller and much more certain. **D** Comparison between the inferred clock rates. The coalescent with reassortment infers lower clock rates.

213 **Contrasting reassortment rates across different hu-** 214 **man influenza viruses**

215 We compared the reassortment rates of different influenza types. To do
216 so, we used the same datasets as described above, as well as an influenza
217 A/H2N2 dataset sampled between 1957 and 1970. We then jointly inferred
218 the reassortment network, the embedding of segment trees, evolutionary
219 rates, effective population sizes and reassortment rates of these viruses. We
220 find that the estimated reassortment rates vary greatly between different
221 influenza viruses.

222 Influenza A/H3N2 shows the highest rates of reassortment, while pan-
223 demic 1918 like H1N1 and influenza B show the lowest inferred rates of
224 reassortment (see figure 3). H2N2 and 2009 pandemic like H1N1 show
225 intermediate rates of reassortment, although the uncertainty on those es-
226 timates is quite large.

227 Differences between p1918-like H1N1, H2N2 and H3N2 are particu-
228 larly interesting since these strains share many common segments. All
229 segments with the exception of HA, NA, and PB1 of influenza A/H3N2
230 originate from the p1918-like H1N1 strain (Scholtissek et al., 1978) and
231 H2N2 and H3N2 only differ in HA and PB1. Pandemic 2009-like human
232 H1N1 which became seasonal in the years after the 2009 pandemic on
233 the other hand has one segment (PB1) that originates from human H3N2
234 and three segments (HA, NP, and NS) derived from classic swine viruses
235 which are descended from a p1918-like strain (Smith, Vijaykrishna, et al.,
236 2009). It shows similar reassortment rates to H3N2, but highly elevated
237 levels compared to the p1918-like H1N1 strain.

238 Such variations in reassortment rates may be driven by a number of
239 factors. Differences in co-infection rate (which may be linked to the ef-
240 fective population size) lead to different probabilities of viruses being in
241 the same host at the same time and therefore to difference in the rate at
242 which reassortants appear. In particular, the higher incidence of Influenza
243 A/H3N2 and the correspondingly likely higher number of co-infection
244 events compared to other influenza A viruses or influenza B viruses may
245 contribute to the higher observed reassortment rate in that case. Addi-
246 tionally, potentially different survival probabilities of reassortants could
247 affect the observed reassortment rates.

248 **Reassortment events occur on fitter parts of reas-** 249 **sortment networks**

250 Next, we test if there is a fitness effect associated with reassortment events.
251 To do so, we classify every network edge from the posterior distribution of
252 inferred networks as either “fit” or “unfit”. We define a fit edge to be any
253 edge having descendants which still persist at least 2 years into the future,
254 while every other edge in the network is defined to be unfit. If reassortment
255 events are beneficial, lineages that are the result of reassortment events
256 should have a higher survival probability and are therefore more likely to
257 persist further into the future.

258 To test if this is the case, we calculated the number of reassortment
259 events on fit edges and on non-fit edges for all networks in the posterior
260 distribution of the MCMC. We then divided this number by the total
261 length of fit respectively not fit edges. As shown in figure 4A, reassortment
262 events occur at a higher rate on fit edges of the H3N2 and influenza B
263 networks than they do on non-fit edges. This suggests that reassortment

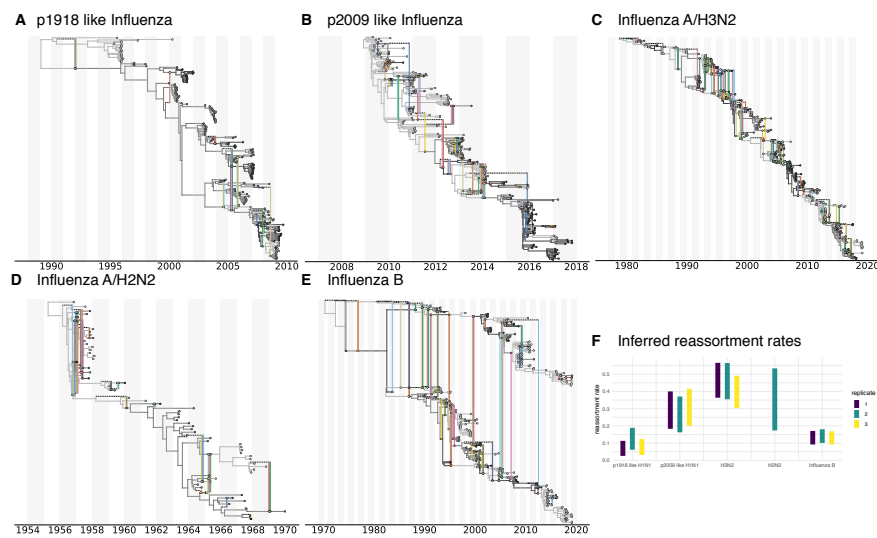


Figure 3: **Estimates of maximum clade credibility networks and reassortment rates of different human influenza viruses.** Maximum clade credibility (mcc) networks of p1918 Influenza A/H1N1 **A**, p2009 influenza A/H1N1 **B**, influenza A/H2N2 **C**, influenza A/H3N2 **D** and influenza B **E**. These mcc networks are shown for one of the random subsets. The mcc networks of all random subsets are shown in figures S3-S6. **F** Here we show the inferred reassortment rates (y-axis) for different influenza viruses on the x-axis. The reassortment rates are per lineage and year.

264 is beneficial to the fitness of influenza A/H3N2 and influenza B viruses.

265 For the other human influenza viruses, fitness benefits of reassortment
266 are less pronounced. For p09-like H1N1 and H2N2, the sampling time
267 windows (both) or number of samples (H2N2) was however rather small
268 and the results are likely driven by a lack of data. Non p09-like H1N1
269 on the other hand had relatively few reassortment events overall driven
270 by a low reassortment rate and the results are likely driven by a lack of
271 reassortment events. These patterns largely hold if the definition of what
272 is a fit edges is changed to having descendants at least 4 or 6 years into
273 the future (see figure S7 & S8). It however decreases for p09 like H1N1
274 for which the overall sampling interval is only 10 years.

275 Since H3N2 has been densely sampled over long time intervals, we
276 analysed two more influenza A/H3N2 datasets, one sampled between 1980
277 and 2005 and one sampled from 2005 until today. For both these datasets,
278 we find higher rates of reassortment on fit edges (see figure S9) . We next
279 tested if for datasets sampled over short times (2 years), we would estimate
280 reassortment rates consistent with the estimated rates on unfit edges. To
281 do so, we compiled 9 dataset, each with 100 to 200 sequences sampled from
282 2 seasons between 2000 and 2018. Averaged over all 9 datasets, we find
283 the short-term reassortment rate to be approximately 0.2 reassortment
284 events per lineage per year, which is consistent with the reassortment rate
285 estimates for unfit edges (see figure S10).

286 Finally, we sought to rule out the possibility that these patterns are
287 simply a property of our reassortment model. To do this, we simulated net-
288 works under the coalescent with reassortment with the reassortment rates
289 and effective population sizes fixed to the mean values estimated from the
290 empirical data, and the network leaf times fixed to those from the same
291 data. We then recomputed the same fit/unfit reassortment rate statis-
292 tics from these simulated networks (see figure S11) and found that the
293 patterns we observed in the empirical data observed patterns completely
294 disappeared. This strongly suggests that the elevated rate of reassortment
295 on fit lineages is not due to the particulars of our model, but is instead a
296 real effect.

297 Conclusion

298 We here present a novel Bayesian approach to jointly infer the reassort-
299 ment network, the embedding of segment trees and the corresponding
300 evolutionary parameters. We show that this approach is able to retrieve
301 reassortment rates, effective population size and reassortment events from
302 simulated data.

303 We have used this facility to show that there are larger differences in
304 the rates of reassortment across different Influenza viruses, and that reas-
305 sortment events occur predominantly in fitter parts of the corresponding
306 reassortment networks. We propose that this is due to selection favour-
307 ing lineages that have reassorted. Although we have deployed a relatively
308 simple way of defining which edges of a network are fit and which are
309 not, future approaches could more directly incorporate fitness models into
310 these network type approaches (Łuksza and Lässig, 2014; Neher, Russell,
311 and Shraiman, 2014).

312 Even if one is not directly interested in reassortment patterns, our
313 approach allows phylogenetic and phylodynamic inferences to exploit full
314 genome sequences for reassorting viruses. This helps to avoid bias and

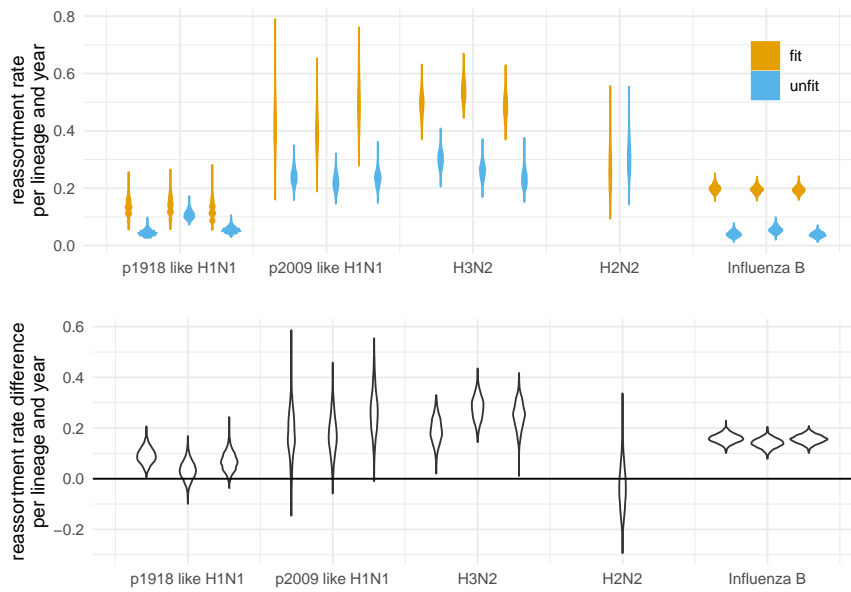


Figure 4: **Estimates of reassortment rates on fit and unfit edges. A** Here we show the number of reassortment events on fit and unfit edges of the networks divided by the total length of fit and unfit edges. Fit edges are defined as having sampled descendant at least 2 years into the future. Every other edge is considered unfit. These rates are shown for different human influenza viruses on the x-axis. The violin plots denote the distribution of these ratios over the posterior distribution of networks. **B** Here we show the difference between fit and unfit reassortment rates. Values above 0 indicate that reassortment events are more likely to occur on fitter, while values below 0 indicate that reassortment events are more likely to occur less fit edges.

315 increase precision compared to, for instance, assuming segments evolve
316 completely independently. However, a lot of development remains to be
317 done in the direction of incorporating skyline models for population size
318 dynamics (Drummond et al., 2005; Minin, Erik W Bloomquist, and Marc
319 A Suchard, 2008) together with extending the model to account for popu-
320 lation structure (Vaughan, Kühnert, et al., 2014; Müller, D. A. Rasmussen,
321 and Stadler, 2017).

322 In summary, this approach allows us to perform network inference
323 by directly accounting for a special kind of recombination process, i.e.
324 reassortment. In the future, we will pursue the development of related
325 approaches to account for a variety of other recombination processes.

326 Methods and Materials

327 The coalescent with reassortment

328 Here we introduce a model to describe a coalescent process with reassort-
329 ment. To do so, we define t to be the time (increasing into the past) before
330 the most recent sample, and L_t as the set of network lineages extant at
331 time t (see Figure 5). Each extant network lineage $l \in L_t$ carries the full
332 set of genome segments, S . In general however, only a subset $\mathcal{C}(l) \subseteq S$ of
333 these are directly ancestral to sampled viruses. We refer to this subset as
334 the *carrying load*. We further define the total number of segments $|S|$ and
335 the number of ancestral segments $|\mathcal{C}(l)|$.

The coalescent with reassortment is a continuous time Markov pro-
cess that proceeds backward in time. It involves three possible events:
sampling, *coalescent* and *reassortment* events. As is usually case for co-
alescent approaches, we condition on sampling events. These happen at
predefined times and simply the number of active network lineages by 1.
Coalescent events occur between two network lineages l and l' at a rate
that is inversely proportional to the effective population size N_e and re-
duce the number of active network lineages by 1. The smaller the effective
population size, the more likely two lineages are to share a common an-
cestor, i.e. the more likely they are to coalesce. Upon a coalescent event,
the segments that the parent lineage p of lineages l and l' carries is the
union of the segments that is carried by l and l' , i.e.:

$$\mathcal{C}(p) = \mathcal{C}(l) \cup \mathcal{C}(l')$$

336 This coalescent events in the network only corresponds to a coalescent
337 event in a segment tree when the corresponding segment is present in
338 both $\mathcal{C}(l)$ and $\mathcal{C}(l')$.

339 Reassortment events happen at a rate ρ per lineage per unit time.
340 A reassortment event on lineage l will increase the number of network
341 lineages by 1. The segments carried by lineage l are randomly assigned to
342 the two parent lineages p_1 and p_2 . This means that the probability of the
343 ancestry of a given segment to follow p_1 , for example, is 0.5.

As we are not interested in the history of segments that are not an-
cestral to our sample, we explicitly integrate over this ancestry in our
model. As with standard coalescent with recombination models, this is
done by omitting non-ancestral events from the process and modifying
the reassortment rate to exactly account for this omission. In our model,
the events which are omitted are “reassortment” events on l in which the
ancestry of every ancestral lineage in $\mathcal{C}(l)$ is assigned to the same parent.
(Thus no true reassortment occurs.) Since each segment chooses its parent

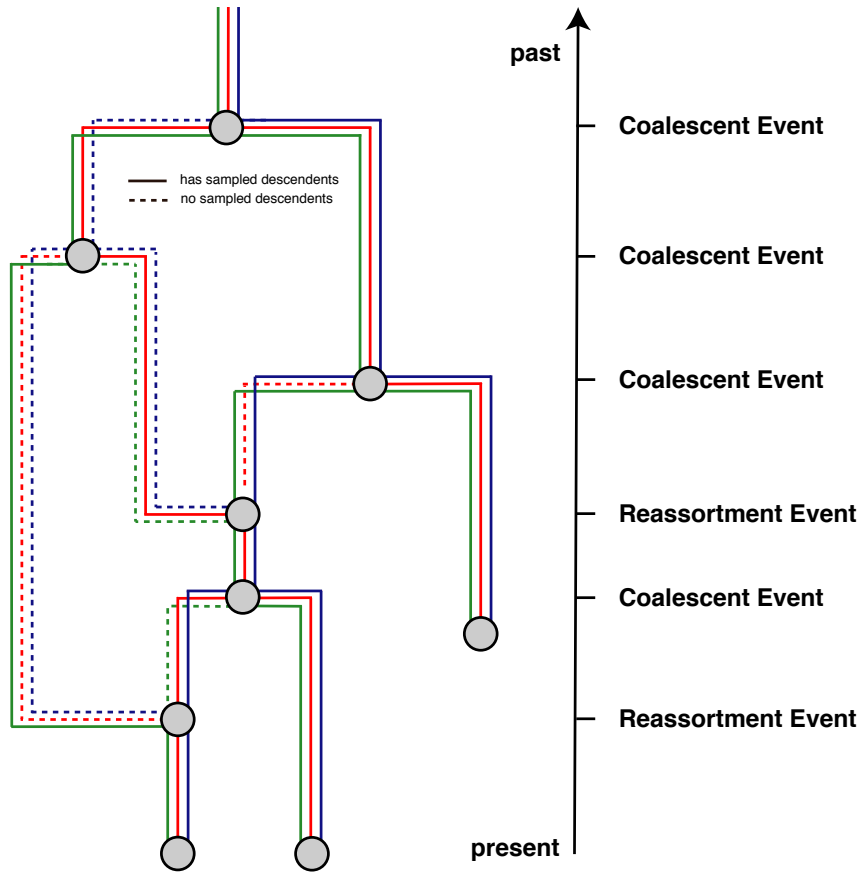


Figure 5: **Example reassortment network.** Here we give an example of a reassortment network where we track 3 different segments differentiated by the different colors through the network. Dashed lines denote segment lineages that do not have sampled descendants. As done in coalescent approaches, we track the network from the present backwards in time to the past.

edge uniformly at random, the probability of either p_1 or p_2 being chosen as ancestral to all segments is

$$P(\mathcal{C}(p_1) = \emptyset \vee \mathcal{C}(p_2) = \emptyset) = 2 \times \left(\frac{1}{2}\right)^{|\mathcal{C}(l)|} = f(l)$$

344 The effective rate of “observable” reassortments on lineage l is then simply
 345 $\rho(1 - f(l))$.

346 Calculating the posterior probability

In order to perform joint Bayesian inference of reassortment networks together with the parameters of the associated models, we use a Markov chain Monte Carlo (MCMC) algorithm to characterize the joint posterior density

$$P(N, \vec{\mu}, \theta, \rho | \vec{D}) = \frac{P(\vec{D} | N, \vec{\mu})P(N | \theta, \rho)P(\vec{\mu}, \theta, \rho)}{P(\vec{D})}. \quad (1)$$

347 Here N represents the full reassortment network (including the embedding
348 of the segment trees), the elements of the vectors \vec{D} and $\vec{\mu}$ represent
349 the segment-specific multiple sequence alignments and their associated
350 molecular substitution models and parameters. The parameters θ and ρ
351 are the effective population size and per-lineage reassortment rate.

352 The terms on the right-hand side of Eq. (1) include the network like-
353 lihood $P(\vec{D}|N, \vec{\mu})$, the network prior $P(N|\rho, \theta)$ and the joint parameter
354 prior $P(\vec{\mu}, \theta, \rho)$. Each of these terms is discussed below. (The denomina-
355 tor $P(\vec{D})$ is the marginal likelihood of the model and does not concern us
356 here.)

357 The network likelihood

358 The usual conditional independence of sites assumption made in phylo-
359 genetic analyses allows us to factorize the network likelihood in terms of
360 the individual segment tree likelihoods:

$$361 P(\vec{D}|N, \vec{\mu}) = \prod_{s \in S} P(D_s | T_s, \mu_s).$$

362 These tree likelihoods can be computed using the standard pruning algo-
363 rithm (Felsenstein, 1981).

364 The network prior

365 The term $P(N|\theta, \rho)$ denotes the probability of the network and the em-
366 bedding of segment trees under the coalescent with reassortment model,
367 with effective population size θ and per-lineage reassortment rate ρ . It
368 plays the role of the tree prior in standard phylodynamic analyses.

369 We can calculate $P(N|\theta, \rho)$ by expressing it as the product of exponen-
370 tial waiting times between reassortment, coalescent and sampling events,
371 i.e.:

$$372 P(N|\theta, \rho) = \prod_{i=1}^{\text{\#events}} P(\text{event}_i | L_i, \theta, \rho) \times P(\text{interval}_i | L_i, \theta, \rho)$$

373 where we define t_i to be the time of the i^{th} event, and L_i to be the set
374 of lineages extant immediately prior to this event. (That is, $L_i = L_t$ for
375 $t \in [t_{i-1}, t_i)$.)

376 **Event contribution.** The event contribution of the i^{th} event in the
377 network is different depending on if the i^{th} event is a coalescent or reas-
378 sortment event. If the i^{th} event is a coalescent event between lineage l_1
379 and l_2 , the event contribution is the probability density of this particu-
380 lar pair of lineages coalescing at time t_i . For a constant-sized coalescent
381 model, this is

$$382 P(\text{event}_i | L_i, \theta, \rho) = \frac{1}{\theta}.$$

383 On the other hand, if the i^{th} event is a reassortment event on lineage l ,
384 the event contribution is the probability density of an (observable) reas-
385 sortment event to occur on that lineage, i.e:

$$386 P(\text{event}_i | L_i, \theta, \rho) = \rho \left[1 - 2 \times \left(\frac{1}{2} \right)^{|C(l)|} \right]$$

387 As we condition on sampling events, their event contribution is always
388 simply 1.

389 **Interval contribution** The interval contribution $P(\text{interval}_i | L_i, \theta, \rho)$
390 is the probability of not observing any event in a given time interval. Three
391 different types of events can happen in the coalescent with reassortment:
392 sampling, coalescent and reassortment events. Since we condition on the
393 times of the sampling events, only coalescent and reassortment events are
394 produced by the CTMC. Given the total rate Λ_i (probability per unity
395 time) with which these occur in the interval immediately prior to event i ,
396 the interval contribution can be written as

$$397 \quad P(\text{interval}_i | L_i, \theta, \rho) = \exp[-\Lambda_i(t_i - t_{i-1})].$$

398 The total rate is the sum of the coalescence rate $\lambda_i^{(c)}$ and the reassortment
399 rate $\lambda_i^{(r)}$. The coalescence rate depends on the number of lineages extant
400 at a particular time and the effective population size in the usual way.

$$401 \quad \lambda_i^{(c)} = \binom{|L_i|}{2} \frac{1}{\theta}.$$

402 The rate of observable reassortment events is

$$403 \quad \lambda_i^{(r)} = \rho \left[|L_i| - \sum_{l \in L_i} \left(\frac{1}{2} \right)^{|C(l)|-1} \right].$$

404 Note that this is generally less than the total rate of reassortment events
405 in this interval, which would be simply $\rho |L_i|$, as this rate excludes reas-
406 sortment events that produce lineages carrying no ancestral segments.

407 **The parameter priors**

408 The term $P(\vec{\mu}, \theta, \rho)$ denotes joint prior distribution of all model param-
409 eters. We factorize this, writing it as the product of the individual param-
410 eter priors $P(\vec{\mu})$, $P(\theta)$ and $P(\rho)$. This asserts that our prior information on
411 any one of these model parameters is independent of the prior information
412 we have for the others.

413 **An MCMC algorithm for reassortment networks**

414 In order to perform MCMC sampling of network and the embedding of
415 segment trees within these networks, we introduce several MCMC oper-
416 ators. These operators often have analogues in operators used to explore
417 different phylogenetic trees. We here only briefly discuss what these oper-
418 ators are continually doing and provide more details in the supplement:

- 419 1. Add/Remove operators which add and remove reassortment events
420 extends the SPR move for networks (Bordewich, Linz, and Semple,
421 2017) to jointly operate on segment trees as well.
- 422 2. Segment diversion operators which change the path segments take
423 at reassortment events.
- 424 3. Exchange operators which change the attachment of edges in the
425 network while keeping the network length constant.
- 426 4. Sub-network Slide operators which change the height of nodes in the
427 network while allowing to change the topology
- 428 5. Scale operators which scale the heights of individual nodes or the
429 whole network without changing the network topology.

- 430 6. Re-simulating above the segment tree roots operator, to efficiently
431 sample parts of the network that are not informed by any genetic
432 data.
- 433 7. Empty segment operator to augment the network with edges that do
434 not carry any segments for the duration of a move, to allow larger
435 jumps in network space.

436 We validate the the implementation of the coalescent with reassortment
437 network prior as well as all operators in the supplement.

438 Summarizing reassortment networks

439 To summarize over a distribution of networks, we use a similar strategy
440 to the maximum clade credibility strategy used to summarize over dis-
441 tributions of trees. To do so, we first compute all unique coalescent and
442 reassortment nodes was encountered during the MCMC. To do so, we have
443 to define when two coalescent or reassortment nodes are the same. We de-
444 fine two coalescent nodes to be the same if a) the parent edges of those
445 nodes carry the same segments and b) if the sub-tree below each segment
446 includes exactly the same clades between the two coalescent nodes. We
447 define two reassortment nodes to be the same if a) both parent edges carry
448 the same segments in the same relative orientation and b) if the sub-tree
449 below each segment includes exactly the same clades between the two
450 reassortment events. This however also means that the more segments
451 we include in the summary, the more likely two nodes will be considered
452 different nodes.

453 While the number of coalescent and reassortment nodes in the network
454 changes over the course of the MCMC, the number of coalescent nodes on
455 the segment trees is constant. In order to avoid dimensionality issues when
456 summarizing, we first compute the frequency of observing each coalescent
457 node over the course of the MCMC. We then weight this frequency by the
458 number of coalescent events on segment trees this coalescent node corre-
459 sponds to. We next choose the network that maximizes those weighted
460 clade credibilities as the maximum clade credibility (or MCC) network.

461 In order to compute the posterior support of each reassortment event
462 in the MCC network, we next compute the frequency of observing each
463 reassortment event in the MCC network during the MCMC.

464 Since we require the network to be rooted, we track segments event
465 after the root of a segment tree was reached. These patterns are however
466 not supported by any genetic information and follow the prior distribu-
467 tion only. For the summary of networks, we therefore remove segments
468 from edges if the root of a segment tree has been reached. Additionally,
469 we remove reassortment loops, i.e. events that start on one edge and then
470 directly reattach to the same edge. Since the support for individual events
471 can greatly depend on how many segments are analysed, we also imple-
472 mented the option to only summarize over a subset of the segments, while
473 ignoring others.

474 Reassortment distance

475 For any reassortment event where segments a and b take different paths,
476 we compute the reassortment distance of segment a onto segment b as fol-
477 lows: First, we follow segment a until it reaches a network edge that carries
478 segment b . We then compute the common ancestor height between seg-
479 ment b at the reassortment event and segment b on that network edge. The

480 reassortment distance of segment a onto segment b is then the different
481 between this common ancestor height and the height of the reassortment
482 event. This seeks to denote for how long segment b in the two parent
483 viruses at the reassortment event evolved independently.

484 Implementation

485 We implemented the MCMC framework for the Coalescent with Reas-
486 sortment as a BEAST2 package called CoalRe. This package includes the
487 classes to do simulation and inference under the coalescent with reassort-
488 ment. The implementation is such that the tree likelihood calculations are
489 separate from the the network framework, which allows to make use of the
490 vast amount of different site and clock models implemented in BEAST2.
491 Additionally, it can be used with other Bayesian approach such as Nested
492 Sampling or coupled MCMC. Further, model comparison as well as inte-
493 gration over evolutionary models can be performed. The package can be
494 downloaded by using the package manager in BEAUti. The source code
495 for the software package can be found here: [https://github.com/nicfel/](https://github.com/nicfel/CoalRe)
496 [CoalRe](https://github.com/nicfel/CoalRe). A tutorial on how to set-up an analysis using the coalescent with
497 reassortment is available at [https://taming-the-beast.org/tutorials/](https://taming-the-beast.org/tutorials/Reassortment-Tutorial/)
498 [Reassortment-Tutorial/](https://taming-the-beast.org/tutorials/Reassortment-Tutorial/) (Barido-Sottani et al., 2017). Networks are
499 logged in the extended Newick format (Cardona, Rosselló, and Valiente,
500 2008) and can be visualized using for example icytree.org (Vaughan,
501 2017). Additionally, we provide python scripts to plot networks based
502 on <https://github.com/evogytis/baltic>.

503 Datasets and data availability

504 We compiled datasets from several influenza viruses using sequence data
505 downloaded from fludb.org (pandemic and seasonal H1N1, H3N2 and in-
506 fluenza B). For the influenza A/H2N2 dataset, we ended up using the same
507 sequences as in (Joseph et al., 2015). We downloaded these sequences from
508 gisaid.org (acknowledgement table can be found here) For all datasets,
509 but the influenza A/H2N2 dataset, we first sub-sampled all sequences to
510 end up with at least 500 samples sampled evenly over time. We then
511 aligned all segments using Muscle 3.8.31 (Edgar, 2004).

512 We then analysed every influenza virus under the coalescent with
513 reassortment in BEAST 2.5.2 (Bouckaert et al., 2019) using coupled
514 MCMC (Altekar et al., 2004; Mueller and Bouckaert, 2019). We assumed
515 the sequences to have evolved under an HKY + Γ_4 model (Hasegawa,
516 Kishino, and Yano, 1985; Yang, 1993), allowing the first two codon posi-
517 tion and the third having different rates (Shapiro, Rambaut, and Drum-
518 mond, 2005). We then jointly estimated all evolutionary rates, the reas-
519 sortment networks and embedding of segments trees, as well as the reas-
520 sortment rates and effective population sizes. For the influenza A/H2N2
521 dataset, we additionally estimated the sampling times for all sequences
522 for which only the year in which the sample was taken was known.

523 For virus types with sequences downloaded from fludb.org, the full
524 XML files to run the datasets are available online. For the influenza
525 A/H2N2 sequences that were obtained from gisaid.org, we remove the
526 sequence characters from the XML files in order to comply with li-
527 cence regulations of gisaid.org. Apart from the sequence characters, the
528 XML files are complete. All other data, such as log files of BEAST2

529 runs, as well as scripts to analyse and plot results are available here
530 <https://github.com/nicfel/Reassortment-Material>.

531 Acknowledgement

532 We would like to thank Alexei J. Drummond, Simone Linz and Daniel
533 Huson for useful discussions on how to summarize networks. NFM and TS
534 are funded by the Swiss National Science foundation (SNF; grant number
535 CR32I3_166258).

536 References

- 537 Altekar, Gautam et al. (2004). “Parallel metropolis coupled Markov
538 chain Monte Carlo for Bayesian phylogenetic inference”. In:
539 *Bioinformatics* 20.3, pp. 407–415.
- 540 Barido-Sottani, Joëlle et al. (2017). “Taming the BEAST? A com-
541 munity teaching material resource for BEAST 2”. In: *Systematic
542 biology* 67.1, pp. 170–174.
- 543 Bloomquist, Erik W. and Marc A. Suchard (2010). “Unifying vertical
544 and nonvertical evolution: a stochastic ARG-based framework.”
545 eng. In: *Syst Biol* 59.1, pp. 27–41. DOI: 10.1093/sysbio/syp076.
546 URL: <http://dx.doi.org/10.1093/sysbio/syp076>.
- 547 Bordewich, Magnus, Simone Linz, and Charles Semple (2017). “Lost
548 in space? Generalising subtree prune and regraft to spaces of
549 phylogenetic networks”. In: *Journal of theoretical biology* 423,
550 pp. 1–12.
- 551 Bouckaert, Remco et al. (2019). “BEAST 2.5: An advanced software
552 platform for Bayesian evolutionary analysis”. In: *PLoS computa-
553 tional biology* 15.4, e1006650.
- 554 Cardona, Gabriel, Francesc Rosselló, and Gabriel Valiente (2008).
555 “Extended Newick: it is time for a standard representation of
556 phylogenetic networks”. In: *BMC bioinformatics* 9.1, p. 532.
- 557 De Vienne, Damien M (2018). “Tanglegrams are misleading for vi-
558 sual evaluation of tree congruence”. In: *Molecular biology and
559 evolution* 36.1, pp. 174–176.
- 560 Didelot, Xavier et al. (2010). “Inference of Homologous Recombina-
561 tion in Bacteria Using Whole-Genome Sequences”. In: *Genetics*
562 186, p. 1435. DOI: 10.1534/genetics.110.120121. URL: <https://dx.doi.org/10.1534/genetics.110.120121>.
- 564 Drummond, Alexei J et al. (2005). “Bayesian coalescent inference of
565 past population dynamics from molecular sequences”. In: *Molec-
566 ular biology and evolution* 22.5, pp. 1185–1192.
- 567 Dudas, Gytis and Trevor Bedford (2019). “The ability of single genes
568 vs full genomes to resolve time and space in outbreak analysis”.
569 In: *bioRxiv*, p. 582957.
- 570 Dudas, Gytis, Trevor Bedford, et al. (2014). “Reassortment between
571 influenza B lineages and the emergence of a coadapted PB1–
572 PB2–HA gene complex”. In: *Molecular biology and evolution*
573 32.1, pp. 162–172.

- 574 Edgar, Robert C (2004). “MUSCLE: multiple sequence alignment
575 with high accuracy and high throughput”. In: *Nucleic acids re-*
576 *search* 32.5, pp. 1792–1797.
- 577 Felsenstein, J. (1981). “Evolutionary trees from DNA sequences:
578 a maximum likelihood approach.” eng. In: *J Mol Evol* 17.6,
579 pp. 368–376. URL: [http://www.ncbi.nlm.nih.gov/pubmed/](http://www.ncbi.nlm.nih.gov/pubmed/7288891)
580 [7288891](http://www.ncbi.nlm.nih.gov/pubmed/7288891).
- 581 Guan, Yi et al. (2010). “The emergence of pandemic influenza
582 viruses”. In: *Protein & cell* 1.1, pp. 9–13.
- 583 Hasegawa, Masami, Hirohisa Kishino, and Taka-aki Yano (1985).
584 “Dating of the human-ape splitting by a molecular clock of
585 mitochondrial DNA”. In: *Journal of molecular evolution* 22.2,
586 pp. 160–174.
- 587 Joseph, Udayan et al. (2015). “Adaptation of pandemic H2N2 in-
588 fluenza A viruses in humans”. In: *Journal of virology* 89.4,
589 pp. 2442–2447.
- 590 Jukes, Thomas H, Charles R Cantor, et al. (1969). “Evolution of pro-
591 tein molecules”. In: *Mammalian protein metabolism* 3.21, p. 132.
- 592 Lu, Lu, Samantha J Lycett, and Andrew J Leigh Brown (2014). “Re-
593 assortment patterns of avian influenza virus internal segments
594 among different subtypes”. In: *BMC evolutionary biology* 14.1,
595 p. 16.
- 596 Luksza, Marta and Michael Lässig (2014). “A predictive fitness
597 model for influenza”. In: *Nature* 507.7490, p. 57.
- 598 McDonald, Sarah M et al. (2016). “Reassortment in segmented RNA
599 viruses: mechanisms and outcomes”. In: *Nature Reviews Micro-*
600 *biology* 14.7, p. 448.
- 601 McVean, Gilean A T and Niall J Cardin (2005). “Approximating the
602 coalescent with recombination.” eng. In: *Philos Trans R Soc Lond*
603 *B Biol Sci* 360.1459, pp. 1387–1393. DOI: 10.1098/rstb.2005.
604 [1673](http://dx.doi.org/10.1098/rstb.2005.1673). URL: <http://dx.doi.org/10.1098/rstb.2005.1673>.
- 605 Minin, Vladimir N, Erik W Bloomquist, and Marc A Suchard (2008).
606 “Smooth skyride through a rough skyline: Bayesian coalescent-
607 based inference of population dynamics”. In: *Molecular biology*
608 *and evolution* 25.7, pp. 1459–1471.
- 609 Mueller, Nicola Felix and Remco Bouckaert (2019). “Coupled
610 MCMC in BEAST 2”. In: *bioRxiv*, p. 603514.
- 611 Müller, Nicola F, David A Rasmussen, and Tanja Stadler (2017).
612 “The structured coalescent and its approximations”. In: *Molec-*
613 *ular biology and evolution* 34.11, pp. 2970–2981.
- 614 Nakajima, Katsuhisa, Ulrich Desselberger, and Peter Palese (1978).
615 “Recent human influenza A (H1N1) viruses are closely related ge-
616 netically to strains isolated in 1950”. In: *Nature* 274.5669, p. 334.
- 617 Neher, Richard A, Colin A Russell, and Boris I Shraiman (2014).
618 “Predicting evolution from the shape of genealogical trees”. In:
619 *Elife* 3, e03568.
- 620 Nelson, Martha I et al. (2008). “Multiple reassortment events in the
621 evolutionary history of H1N1 influenza A virus since 1918”. In:
622 *PLoS pathogens* 4.2, e1000012.

- 623 Rasmussen, Matthew D et al. (2014). “Genome-wide inference of an-
624 cestral recombination graphs”. In: *PLoS genetics* 10.5, e1004342.
- 625 Scholtissek, C et al. (1978). “On the origin of the human influenza
626 virus subtypes H2N2 and H3N2”. In: *Virology* 87.1, pp. 13–20.
- 627 Shapiro, Beth, Andrew Rambaut, and Alexei J Drummond (2005).
628 “Choosing appropriate substitution models for the phylogenetic
629 analysis of protein-coding sequences”. In: *Molecular biology and
630 evolution* 23.1, pp. 7–9.
- 631 Smith, Gavin JD, Justin Bahl, et al. (2009). “Dating the emergence
632 of pandemic influenza viruses”. In: *Proceedings of the National
633 Academy of Sciences* 106.28, pp. 11709–11712.
- 634 Smith, Gavin JD, Dhanasekaran Vijaykrishna, et al. (2009). “Ori-
635 gins and evolutionary genomics of the 2009 swine-origin H1N1
636 influenza A epidemic”. In: *Nature* 459.7250, p. 1122.
- 637 Steel, John and Anice C Lowen (2014). “Influenza A virus reassort-
638 ment.” In: *Current topics in microbiology and immunology* 385,
639 pp. 377–401. ISSN: 0070-217X. DOI: 10.1007/82_2014_395.
- 640 Vaughan, Timothy G (2017). “IcyTree: rapid browser-based visual-
641 ization for phylogenetic trees and networks”. In: *Bioinformatics*
642 33.15, pp. 2392–2394.
- 643 Vaughan, Timothy G, Denise Kühnert, et al. (2014). “Efficient
644 Bayesian inference under the structured coalescent”. In: *Bioin-
645 formatics* 30.16, pp. 2272–2279.
- 646 Vaughan, Timothy G, David Welch, et al. (2017). “Inferring Ancestral
647 Recombination Graphs from Bacterial Genomic Data”. In:
648 *Genetics* 205 (2), pp. 857–870. ISSN: 1943-2631. DOI: 10.1534/
649 *genetics*.116.193425.
- 650 Westgeest, Kim B et al. (2014). “Genomewide analysis of reassort-
651 ment and evolution of human influenza A (H3N2) viruses cir-
652 culating between 1968 and 2011”. In: *Journal of virology* 88.5,
653 pp. 2844–2857.
- 654 Yang, Ziheng (1993). “Maximum-likelihood estimation of phylogeny
655 from DNA sequences when substitution rates differ over sites.”
656 In: *Molecular biology and evolution* 10.6, pp. 1396–1401.