

1 On the importance of skewed offspring distributions and
2 background selection in viral population genetics

3

4

5 Kristen K. Irwin^{1,2}, Stefan Laurent^{1,2}, Sebastian Matuszewski^{1,2}, Séverine
6 Vuilleumier^{1,2}, Louise Ormond^{1,2}, Hyunjin Shim^{1,2}, Claudia Bank^{1,2,3}, and
7 Jeffrey D. Jensen^{1,2,4}

8

9

10 ¹ – École Polytechnique Fédérale de Lausanne (EPFL), School of Life
11 Sciences, Lausanne, Switzerland

12 ² – Swiss Institute of Bioinformatics (SIB), Lausanne, Switzerland

13 ³ – Instituto Gulbenkian de Ciência (IGC), Oeiras, Portugal

14 ⁴ – Arizona State University (ASU), School of Life Sciences, Center for
15 Evolution & Medicine, Tempe, USA

16

17 Word Count: 3480

18

19 Keywords: Virus, Background Selection, Multiple Merger Coalescent, Skewed
20 Offspring Distribution

21 **Abstract**

22

23 Many features of virus populations make them excellent candidates for
24 population genetic study, including a very high rate of mutation, high levels of
25 nucleotide diversity, exceptionally large census population sizes, and frequent
26 positive selection. However, these attributes also mean that special care must
27 be taken in population genetic inference. For example, highly skewed
28 offspring distributions, frequent and severe population bottleneck events
29 associated with infection and compartmentalization, and strong purifying
30 selection all affect the distribution of genetic variation but are often not taken
31 into account. Here, we draw particular attention to multiple-merger coalescent
32 events and background selection, discuss potential mis-inference associated
33 with these processes, and highlight potential avenues for better incorporating
34 them in to future population genetic analyses.

35

36

37 **Introduction**

38

39 Viruses appear to be excellent candidates for studying evolution in real time;
40 they have short generation times, high levels of diversity often driven by very
41 large mutation rates and population sizes (both census and effective), and
42 they experience frequent positive selection in response to host immunity or
43 antiviral treatment. However, despite these desired attributes, standard
44 population genetic models must be used with caution when making
45 evolutionary inference.

46

47 Firstly, population genetic inference is usually based on a coalescence model
48 of the Kingman type, under the assumption of Poisson-shaped offspring
49 distributions where the variance equals the mean and is always small relative
50 to the population size; consequently, only two lineages may coalesce at a
51 time. In contrast, viruses have highly variable reproductive rates, taken as
52 rates of replication; these may vary based on cell or tissue type, level of
53 cellular differentiation, or stage in the lytic/lysogenic cycle (Knipe and Howley,
54 2007), resulting in highly skewed offspring distributions. This model violation
55 is further intensified by the strong bottlenecks associated with infection and by
56 strong positive selection (Neher and Hallatschek, 2013). Therefore, virus
57 genealogies may be best characterized by *multiple merger* coalescent (MMC)
58 models (e.g, Pitman, 1999; Sagitov, 1999; Donnelly and Kurtz, 1999;
59 Schweinsberg, 2000; Möhle and Sagitov, 2001; Eldon and Wakeley, 2008),
60 instead of the Kingman coalescent.

61

62 Secondly, the mutation rates of many viruses, particularly RNA viruses, are
63 among the highest observed across taxa (Lauring *et al.*, 2013; Cuevas *et al.*,
64 2015). Though these high rates of mutation are what enable new beneficial
65 mutations to arise, potentially allowing for rapid resistance to host immunity or
66 antiviral drugs, they also render high mutational loads (Sanjuán, 2010; Lauring
67 *et al.*, 2013). Specifically, the distribution of fitness effects (DFE) has now
68 been described across taxa – demonstrating that the input of deleterious

69 mutations far outnumbers the input of beneficial mutations (Acevedo *et al.*,
70 2014; Bank *et al.*, 2014; Bernet and Elena, 2015; Jiang *et al.*, 2016). The
71 purging of these deleterious mutants through purifying selection can affect
72 other areas in the genome through a process known as background selection
73 (BGS) (Charlesworth *et al.*, 1993). Accounting for these effects is important for
74 accurate evolutionary inference in general (Ewing and Jensen, 2016), but
75 essential for the study of viruses due to their particularly high rates of mutation
76 and compact genomes (Renzette *et al.*, 2016).

77

78 Given these distinctive features of virus populations and the increasing use of
79 population genetic inference in this area (*e.g.*, Renzette *et al.*, 2013; Foll *et al.*,
80 2014; Pennings *et al.*, 2014; Renzette *et al.*, 2016), it is crucial to account for
81 these processes that are shaping the amount and distribution of variation
82 across their genomes. We aim here to draw particular attention to multiple-
83 merger coalescent events and background selection, and the repercussions
84 of ignoring them in population genetic inference, highlighting particular
85 applications to viruses. We conclude with general recommendations for how
86 best to address these topics in the future.

87

88 ***Skewed Offspring Distributions and the Multiple Merger Coalescent***

89

90 *Inferring evolutionary history using the Wright-Fisher model: benefits and*
91 *shortcomings*

92

93 Many population genetic statistics and subsequent inference are based on the
94 Kingman coalescent and the Wright-Fisher (WF) model (Wright, 1931;
95 Kingman, 1982). With increasing computational power, the WF model has
96 also been implemented in forward-time methods, which allows for the
97 modeling of more complex evolutionary scenarios versus backward-time
98 methods. This also allows for the inference of population genetic parameters,
99 including selection coefficients and effective population sizes (N_e), even from
100 time-sampled data (i.e., data collected at successive time points) (Ewens,
101 1979; Williamson and Slatkin, 1999; Malaspinas *et al.*, 2012; Foll *et al.*, 2014;
102 Foll *et al.*, 2015; Ferrer-Admetlla *et al.*, 2016; Malaspinas, 2016). These
103 methods are robust to some violations of WF model assumptions, such as
104 constant population size, random mating, and non-overlapping generations,
105 and also have been extended to accommodate selection, migration and
106 population structure (Neuhauser and Krone, 1997; Nordborg, 1997; Wilkinson-
107 Herbots, 1998).

108

109 However, it has been suggested that violations of the assumption of a small
110 variance in offspring number in the WF model, and in other models that result
111 in the Kingman coalescent in the limit of large population size, lead to
112 erroneous inference of population genetic parameters (Eldon and Wakeley,
113 2006). Biological factors such as sweepstake reproductive events, population
114 bottlenecks, and recurrent positive selection may lead to skewed distributions
115 in offspring number (Eldon and Wakeley, 2006; Li *et al.*, 2014); examples
116 include various prokaryotes (plague), fungi (*Z. tritici*, *P. striiformis*, rusts,
117 mildew, oomycetes), plants (*A. thaliana*), marine organisms (sardines, cods,
118 salmon, oysters), crustaceans (*Daphnia*), and insects (aphids) (reviewed in
119 Tellier and Lemaire, 2014). The resulting skewed offspring distributions can

120 also result in elevated linkage disequilibrium (LD) despite frequent
121 recombination, as linkage depends not only on recombination rate, but also
122 on the degree of skewness in offspring distributions (Eldon and Wakeley,
123 2008; Birkner *et al.*, 2013). Such events may also skew estimates of F_{ST}
124 relative to those expected under WF models, as there is a high probability of
125 alleles being *identical by descent* in subpopulations, where the expectation of
126 coalescent times within subpopulations is less than that between
127 subpopulations regardless of the timescale or magnitude of gene flow (Eldon
128 and Wakeley, 2009).

129
130 The assumption of small variance in offspring number may often be violated in
131 virus populations as well. For example, progeny RNA virus particles from
132 infected cells can vary up to 100 fold (Zhu *et al.*, 2009). Second, features such
133 as diploidy, recombination, and latent stages are expected to increase the
134 probability of multiple merger events (Davies *et al.*, 2007; Taylor and Véber,
135 2009; Birkner *et al.*, 2013). Third, within their life cycle, viruses experience
136 bottleneck events during transmission and compartmentalization, followed by
137 strong selective pressure from both the immune system and drug treatments.
138 Finally, at the epidemic level, extinction-colonization dynamics drive
139 population expansion (Anderson and May, 1991).

140
141 All of these aspects characterize HIV for example, a diploid virus with
142 extraordinary rates of recombination (Schlub *et al.*, 2014). Transmitted and
143 founder viruses undergo at least two distinct genetic bottlenecks (one of
144 physical transmission and one of infection, respectively; Joseph and
145 Swanstrom, 2015), followed by strong selection imposed by the immune
146 system (Moore *et al.*, 2002). At the epidemic scale, besides multiple events of
147 colonization (Tebit and Arts, 2011), strong heterogeneity in the virus
148 transmission chain has also been observed (*e.g.*, Service and Blower, 1995).

149
150
151

152 *Beyond WF assumptions: the Multiple Merger Coalescent*

153

154 A more general coalescent class of models, summarized as the MMC class,
155 can account for these violations, particularly for (non-Poisson) skewed
156 offspring distributions, by allowing more than two lineages to coalesce at a
157 time (Table 1). These are often derived from Moran models, (Moran, 1958),
158 generalized to allow multiple offspring per individual. In contrast to the
159 Kingman coalescent (for which $P(k > 2) = 0$, where k is the number of
160 lineages coalescing simultaneously), a probability distribution for k -merger
161 events determines coalescence.

162

163 The parameters inferred under the MMC differ from those inferred under the
164 Kingman coalescent in several notable respects. In a Kingman coalescent,
165 effective size N_e scales linearly with census size N , whereas for the MMC it
166 does not (Huillet and Möhle, 2011). Thus genetic diversity is a non-linear
167 function of population size. Coalescent trees under the MMC also have more
168 pronounced star-like genealogies with longer branches (Figure 1), and their
169 site frequency spectra (SFSs) are skewed toward an excess of low frequency
170 and high frequency variants because of these branch lengths (Eldon and
171 Wakeley, 2006; Blath *et al.*, 2016), generating a more negative Tajima's D
172 (Birkner *et al.*, 2013). With similar migration and population size, alleles fix at
173 a higher rate per population in the MMC than under the Kingman coalescent,
174 and thus higher F_{ST} is expected between subpopulations (Eldon and Wakeley,
175 2009). Further, the efficacy of selection increases, as selection acts almost
176 deterministically between multiple merger events; in the Wright-Fisher model,
177 genetic drift counteracts selection fairly strongly (Der *et al.* 2011), but in
178 generalized models where offspring distributions are wide, beneficial
179 mutations may be more likely to escape stochastic loss and thus continue to
180 fixation. Furthermore, the fixation probability of a new mutant with a positive
181 selection coefficient approaches 1 as the population size increases, in stark
182 contrast with traditional expectations under the standard Wright-Fisher model
183 (Der *et al.*, 2011).

184

185 Not accounting for skewed offspring distributions can lead to mis-inference.

186 For instance, Eldon and Wakeley (2006) showed that for Pacific oysters,
187 which have been argued to undergo sweepstake-like reproductive events
188 (Hedgecock, 1994a), the estimated population-wide mutation rate θ inferred
189 under the Kingman coalescent is two orders of magnitude larger than that
190 obtained from the ψ -coalescent (see below) - 9 vs 0.0308, respectively - and,
191 indeed, provides a poor fit to the data.

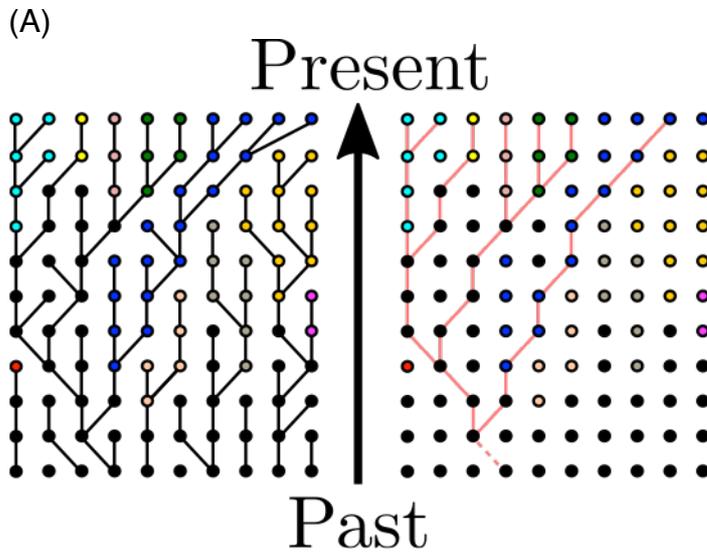
192

193

194 Figure 1: Multiple-Merger and Kingman Coalescent Realizations

195

196

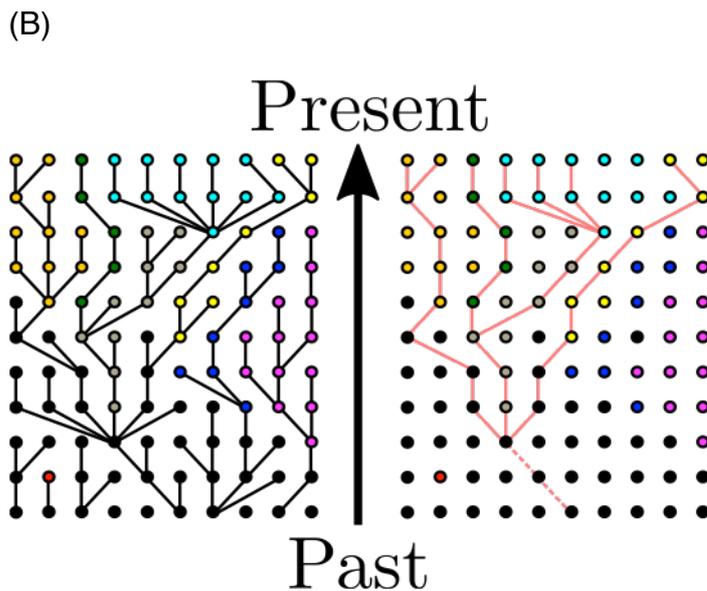


197

198

199

200



201

202

203

204 Figure 1: Example genealogies and samples from (A) the Kingman coalescent

205 and (B) a multiple-merger coalescent. Panels on the left show the evolutionary

206 process of the whole population, whereas those on the right show a possible

207 sampling and its resulting genealogy. Colors correspond to different (neutral)

208 derived allelic states, where black denotes the wild type.

209

210

211 Table 1: Hierarchy of coalescent models, in decreasing order of generality

212

Coalescent model	Allows MMs?	Allows simultaneous MMs?	Distribution and parameters	References
Ξ -coalescent	Yes	Yes	MMC events occur at rate λ , with a specific measure Ξ on the infinite simplex and which allows an arbitrary number of simultaneous mergers.	Schweinsberg (2000); Möhle and Sagitov (2001)
↳ Λ -coalescent	Yes	No	MMC events occur at rate λ (but ≤ 1 event/time)	Donnelly and Kurtz (1999); Pitman (1999); Sagitov (1999)
↳ ψ -coalescent	Yes	No	λ follows a distribution which depends on ψ , i.e., the fraction of the population replaced by the offspring of a single individual	Eldon and Wakeley (2006); Eldon and Wakeley (2008); Eldon and Wakeley (2009); Eldon and Degnan (2012)
↳ Beta-coalescent	Yes	No	λ follows Beta-distribution: $\text{beta}(\alpha, 2-\alpha)$ with $1 \leq \alpha < 2$	Schweinsberg (2003); Berestycki <i>et al.</i> (2007); Berestycki <i>et al.</i> (2008); Birkner and Blath (2008); Birkner <i>et al.</i> (2013); Steinrücken <i>et al.</i> (2013)
↳ Bolthausen-Sznitman	Yes	No	λ follows Beta-distribution with $\alpha=1$: $\text{beta}(1, 1) =$ uniform on $[0, 1]$	Bolthausen and Sznitman (1998); Basdevant and Goldschmidt (2008); Neher and Hallatschek (2013)
↳ Kingman coalescent	No	No	λ follows Beta-distribution with $\alpha=2$; Λ has unit mass at 0 ($\Lambda(dx) = \delta_0(x)dx$)	Kingman (1982)

213 Table 1: Coalescent models listed in decreasing order with respect to
 214 generality; arrows indicate coalescents that are considered subtypes of those
 215 above.

216 *The ψ -coalescent*

217

218 Introduced by Eldon and Wakeley (2006), the ψ -coalescent (also called the
219 'Dirac-coalescent') differentiates two possible reproductive events in the
220 underlying forward process (Figure 2). Either a standard Moran model
221 reproduction event occurs (with probability $1-\varepsilon$), where a single individual is
222 randomly chosen to reproduce and the (single) offspring replaces one
223 randomly chosen non-parental individual; all other individuals, including the
224 parent, persist. Alternatively, a 'sweepstake' reproductive event occurs (with
225 probability ε) (Hedgecock, 1994b), where a single parent replaces $\psi*N$
226 individuals. If these sweepstake events happen frequently enough, the rate of
227 $\psi*N$ -reproduction events will be much greater than that of 2-reproduction
228 events, and the underlying coalescent process will consequently be
229 characterized by MM events; if two or more parents were to replace $\psi*N$
230 individuals, simultaneous MM events may occur in a single generation
231 resulting in a Ξ -coalescent. However, in contrast to other MMC models (*e.g.*,
232 Ξ -coalescent or other λ -coalescents), the parameter ψ has a clear biological
233 interpretation as the fraction of the population that is replaced in each
234 sweepstake reproductive event. Though the assumption of a fixed ψ (as in the
235 normal ψ -coalescent) seems biologically unrealistic, it can be avoided by
236 treating ψ as a Poisson parameter. Finally, despite its appealing connection to
237 biologically relevant measures, the appropriateness of making inferences
238 based on the ψ -coalescent still depends on the biology of the specific virus
239 being studied. Thus, model choice is still essential, and the best-fit coalescent
240 should be assessed on a case-by-case basis.

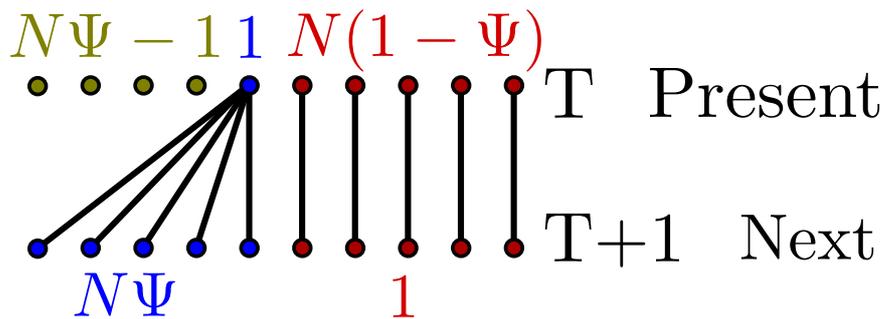
241

242

243

244 Figure 2: Depiction of the modified Moran model underlying the ψ coalescent

245



246

247

248 Figure 2: Lineages between the present and the next generation, where N is
 249 the population size, ϵ is the probability of a sweepstake event, and ψ is the
 250 fraction of the population that is replaced in each such event. Labels in the top row
 251 give the number of parental individuals reproducing in a given manner
 252 (represented by color), whereas labels in the bottom row give the number of
 253 corresponding offspring per parent.

254

255

256 *Application to Viruses*

257

258 There are several reasons why a modified Moran model may better capture
 259 virus evolution than models converging to the Kingman coalescent, although it
 260 does not account for fitness differences between individuals. First, virus
 261 evolution is driven by strong bottlenecks during host transmission and
 262 intrahost selection processes, which likely result in skewed offspring
 263 distributions (Figure 3) (Gutiérrez *et al.*, 2012; Tellier and Lemaire, 2014).
 264 Further, viruses display the MMC-typical low N_e/N ratio (Pennings *et al.*, 2014;
 265 Tellier and Lemaire, 2014), can adapt rapidly (Neher and Hallatschek, 2013),
 266 and may have sweepstake-like reproductive events in which a single virion
 267 can propagate a large fraction of the entire population (Grenfell *et al.*, 2004;
 268 Pybus and Rambaut, 2009). For example, the influenza virus hemagglutinin
 269 (HA) segment appears to be under strong directional selection imposed by
 270 host immunity (and sometimes drug treatment), resulting in a ladder-like

271 genealogy, (as depicted in Figure 3A), suggesting that only a few viruses
272 seed the entire next generation (Grenfell *et al*, 2004). That being said, some
273 challenges remain, such as rigorously defining the term ‘generation’ for virus
274 populations, and subsequently confirming that the per generation mutation
275 rate is on the order of the coalescent timescale c_N , which is a prerequisite for
276 the use of any coalescent approach. Finally, viruses with little or no
277 recombination may be prone to clonal interference, which should be explicitly
278 accounted for in population models and resulting coalescents (*e.g.*, Strelkowa
279 and Lässig, 2012).

280

281 Those processes that make viruses ideal candidates for MMCs can differ by
282 scale (see Figure 3); for example, following transmission events, there are
283 severe founder events and potentially high recombination within the host (*e.g.*,
284 HIV, HCMV). Subsequent compartmentalization may introduce intra-host
285 population structure through bottlenecks, colonization events, and extinction
286 events (Renzette *et al.*, 2013). To date, it remains unclear how often MMCs fit
287 the patterns of variation observed in intra-host relative to inter-host virus
288 populations – but such comparisons are increasingly feasible. Finally, periods
289 of latency - temporary virus inactivation with cessation of reproduction -
290 should be incorporated in such modeling, potentially as recurring mass
291 extinction events (Taylor and Véber, 2009). Thus, multiple MMC models are a
292 necessary but not final step towards addressing the various patterns observed
293 at different scales of virus evolution (Table 1).

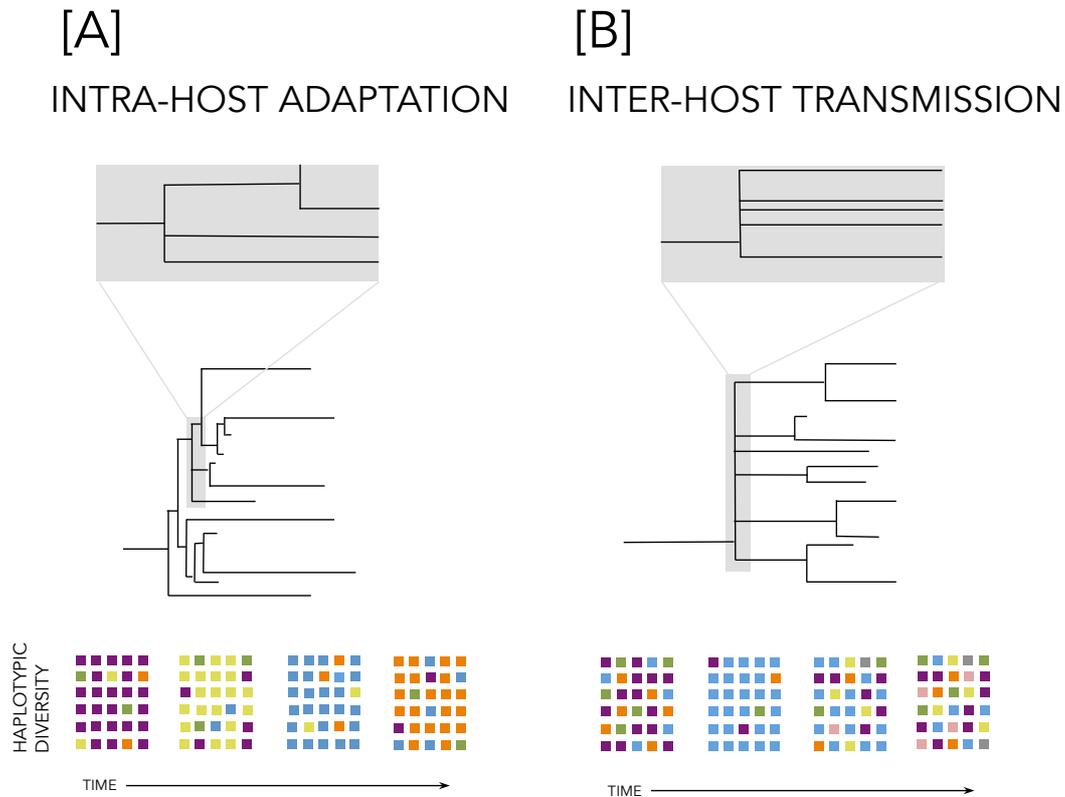
294

295 The large data sets often generated from viruses may also prove impractical
296 for the likelihood-based methods commonly employed for MMCs. This
297 limitation has partially been overcome by Eldon *et al.* (2015), who proposed
298 an approximate likelihood method along with an Approximate Bayesian
299 Computation (ABC) approach based on the SFS to distinguish between the
300 MMC and exponential population growth. Although both effects are expected
301 to result in very similar SFSs, characterized by an excess of singletons as
302 compared to the Kingman coalescent, the bulk and tail of the SFS (*i.e.*, the

303 higher-order frequency classes) typically differ, which can be assessed by
304 approximate likelihood-ratio tests and Approximate Bayes Factors (Eldon *et*
305 *al.*, 2015).

306

307 Figure 3: Example Processes Spurring MM Events in Virus Populations



308

309 Figure 3: Examples include (A) intra-host adaptation (a selective process) and
310 (B) inter-host transmission (a demographic process). The tree in (A)
311 characterizes, for example, NA or HA evolution in the influenza A virus, driven
312 by positive selection; selection by host immunity is ongoing, while that from
313 drug treatment may be intermittent. The tree in (B) represents inter-host
314 transmission and its associated bottleneck; for viruses that compartmentalize
315 (such as HCMV and HIV), similar patterns follow transmission to new
316 compartments. The colored squares below the trees roughly indicate the
317 diversity of the population through time. Intra-host adaptation may temporarily
318 decrease diversity owing to genetic hitchhiking, though single snapshots may
319 not reflect varying temporal levels of diversity. During inter-host transmission,
320 diversity decreases owing to the associated bottleneck but then may quickly
321 recover in the new host.

322

323 [BOX 1: Future challenges in MMC models]

324

325 In order to make MMC models biologically relevant for viruses, a number of
326 important tasks remain:

327

- 328 1. Describe summary statistics that capture demographic features and
329 processes when offspring distributions are highly skewed; such
330 patterns will be required for large-scale inference in a computationally
331 efficient (*e.g.*, Approximate Bayesian) framework.
- 332 2. Better understand the behavior of commonly used summary statistics
333 under such models, as done for F_{ST} by Eldon and Wakeley (2009), for
334 commonly used divergence, SFS, and LD-based statistics.
- 335 3. Determine which MMCs are best suited for different scales of virus
336 evolution (*i.e.*, intra-host, inter-host, global); develop novel models if
337 necessary.
- 338 4. Investigate the effect of violations of MMC assumptions (*e.g.*,
339 overlapping generations, number of multiple merger events) on
340 inference.

341

342 [END BOX 1]

343

344 ***Purifying Selection and Linkage in Virus Populations***

345

346 *Modeling Background Selection*

347

348 The joint modeling of the effects of genetic drift and positive selection,
349 including in experimental evolution studies of virus populations, has improved
350 our ability to distinguish adaptive from neutral mutations by minimizing the
351 chance that the rapid fixation of a neutral allele is incorrectly interpreted as
352 strong positive selection (Li *et al.*, 2012; Foll *et al.*, 2014). However, there is
353 another process that must be incorporated if we are to fully understand
354 mutation trajectories in virus populations: background selection (BGS).

355

356 BGS was originally proposed to explain patterns of reduced diversity in
357 regions of low recombination – patterns that were previously suggested to be
358 the signature of genetic hitchhiking (HH) around strongly beneficial mutations
359 (see Begun and Aquadro, 1992 and Charlesworth *et al.*, 1993). It was argued
360 that only neutral mutations present on the “least-loaded” chromosomes – that
361 is, those with the fewest deleterious mutations – have appreciable
362 probabilities of reaching high frequencies or fixation. Kimura and Maruyama
363 (1966) showed that the proportion of chromosomes belonging to the least-
364 loaded class is

365

$$366 \quad f_0 = \exp\left(-\frac{U}{2hs}\right), \quad (1)$$

367

368 where U is the rate of mutation to a deleterious state, s is the selection
369 coefficient against homozygous mutations, and h is the dominance coefficient.
370 For simplicity of modeling, h is usually set to 1 for viruses that carry a single
371 copy of their genome in each virion, although polyploid effects could arise in
372 the case of multiple virions infecting the same cell.

373

374 The least-loaded class, and thus genetic diversity in the presence of BGS, is
375 dependent on the balance between the influx of deleterious mutations

376 (occurring at rate U) and their removal by natural selection (according to the
377 product hs). Assuming that offspring exclusively originate from the least-
378 loaded class of individuals, Charlesworth *et al.* (1993) expressed the expected
379 neutral diversity due to background selection as

380

$$381 \quad \pi = 4 f_o N_e \mu , \quad (2)$$

382

383 where N_e is the effective population size and μ is the mutation rate. As BGS
384 reduces the number of reproducing individuals, genetic drift increases, thus
385 reducing genetic diversity and increasing stochasticity in allele trajectories.
386 Further, since only the genetic diversity segregating in the least-loaded class
387 can be observed, population size inferred from measures of genetic diversity
388 may be underestimated if BGS is not properly taken into account (Ewing and
389 Jensen, 2016).

390

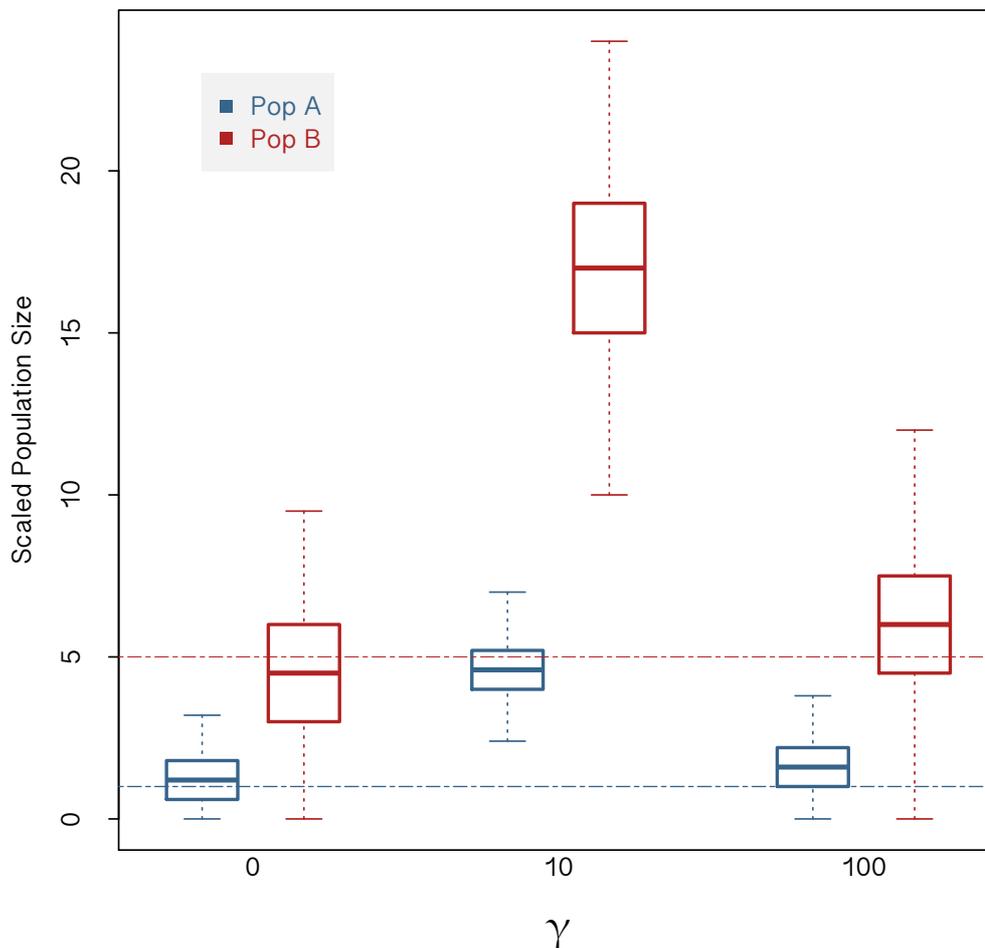
391 In the BGS model described above, strongly deleterious mutations are
392 maintained in mutation-selection balance such that no skew in the SFS is
393 expected, as rare variants are rapidly purged. Thus, a simple re-scaling of N_e
394 is often used as a proxy for the effects of BGS (*e.g.*, Hudson and Kaplan,
395 1995; Zeng and Charlesworth, 2011; Prüfer *et al.*, 2012; Zeng, 2013).

396 However, recent work has demonstrated that, while this re-scaling is
397 appropriate for strongly deleterious mutations, it is largely inappropriate for
398 weakly deleterious mutations that may segregate in the population. Figure 4
399 shows the skew in estimates of population size and migration rates obtained
400 using an ABC approach when BGS is prevalent for two populations A and B
401 that have split at time $\tau=2N_e$ generations (reproduced from Ewing and Jensen,
402 2016). Further, experimental work on the shape of the distribution of fitness
403 effects (DFE) in many organisms indicates that weakly deleterious mutations
404 represent an important class (*e.g.*, Eyre-Walker and Keightley, 2007; Bank *et*
405 *al.*, 2014). These mutations may act to skew the SFS towards rare alleles as
406 they decrease the expected frequency of linked neutral mutations relative to
407 neutral expectations. As subsequent demographic inference is based on the

408 shape of this SFS, this effect should be properly accounted for by directly
409 simulating weakly deleterious mutations rather than implementing a simple
410 rescaling, as is common practice. Though important analytical progress has
411 been made in this area (*e.g.*, McVean and Charlesworth, 2000), simulations
412 remain the best option for the non-equilibrium demographic models and
413 alternative coalescents recommended here for inference in virus populations.

414

415 Figure 4: Bias in parameter inference at intermediate levels of BGS



416

417 Figure 4: Bias in parameter inference for different levels of BGS, redrawn from
418 Ewing & Jensen (2016). Posterior densities from ABC inference for population
419 size are shown. The strength of purifying selection is given as γ , where $\gamma =$
420 $2N_e s$. Population A has a true scaled size of 1 (blue line), and population B a
421 true scaled size of 5 (red line). Both population sizes are scaled relative to the
422 size of the ancestral population. As shown, the greatest mis-inference occurs

423 in the presence of weakly deleterious mutations and subsequent strong BGS
424 effects.

425

426

427

428 *The Effects of Background Selection on Inference in Virus Populations*

429

430 Efforts to estimate the impact of BGS in non-viral organisms have been well
431 reported. One of the most notable examples is that of Comeron (2014), who
432 estimated levels of BGS in *Drosophila melanogaster* based on the results of
433 Hudson and Kaplan (1995) and Nordborg *et al.* (1996) using a high-definition
434 recombination map, with results indicating strong effects across the genome.
435 For viruses, similar efforts are in their infancy, with the first attempt at such
436 estimation in a virus reported recently by Renzette *et al.* (2016), utilizing the
437 theoretical predictions of Innan and Stephan (2003). Interestingly, the full
438 spectrum of recombination frequencies is available in viruses – from non-
439 recombining (*e.g.*, most negative-sense RNA viruses), to re-assorting (*e.g.*,
440 Influenza virus), to rarely recombining (*e.g.*, Hepatitis C and West Nile
441 viruses), to frequently recombining (*e.g.*, HIV), offering a highly promising
442 framework for comparative analyses investigating the pervasiveness of BGS
443 effects (Chare *et al.*, 2003; Simon-Loriere and Holmes, 2011). Further, given
444 the high mutation rates and compact genomes of many viruses, evolutionary
445 theory suggests effects at least equal to those seen in *Drosophila*.

446

447 In order to accomplish such inference, improved recombination maps for virus
448 genomes will be important. With such maps in hand, and given the
449 amenability of viruses to experimental perturbation, it may indeed be feasible
450 to understand and account for BGS in models of virus evolution.

451

452 [BOX 2: Future challenges in identifying the effects of BGS]

453

454 As BGS almost certainly impacts inference in virus populations, accounting for
455 its effects is critical. Future challenges include:

456

- 457 1. Account for BGS effects on the SFS by directly simulating weakly
458 deleterious mutations, rather than by rescaling N_e .
- 459 2. Improve recombination maps for virus genomes.
- 460 3. Develop models combining the effects of non-equilibrium demography,
461 positive selection, and BGS, ideally to allow for the joint estimation of
462 all associated parameters.
- 463 4. Extend methods applied to other taxa to virus populations; for example,
464 establishing a baseline of variation for use as a null expectation to
465 estimate BGS levels across the genome, as done for *Drosophila*.

466

467 [END BOX 2]

468

469

470

471

472

473 ***Future Directions***

474

475 Given that skewed offspring distributions and pervasive linked selection are
476 likely important factors influencing the inference of virus population
477 parameters, it is important to note that multiple backward and forward
478 simulation programs have recently been developed which make the modeling
479 of these processes feasible (Hernandez, 2008; Messer, 2013; Thornton,
480 2014; Eldon *et al.*, 2015; Zhu *et al.*, 2015). This will allow researchers to
481 directly simulate from parameter ranges that may be relevant for their
482 population of interest, developing a better intuition for the importance of these
483 processes in shaping the observed genomic diversity. More concretely, the

484 ability to now simulate in a computationally efficient framework opens the
485 possibility of directly implementing ABC inference approaches under these
486 models. Thus, by drawing mutations from a biologically realistic distribution of
487 fitness effects and allowing offspring distributions to appropriately vary, it is
488 now possible to re-implement common demographic estimation or genome
489 scan approaches; these modified approaches would be based on more
490 appropriate null expectations of the shape of the SFS, the extent of linkage
491 disequilibrium, and the degree of population divergence.

492 **Acknowledgements**

493

494 We would like to thank Bjarki Eldon for helpful suggestions during the early
495 stages of this manuscript. This work was funded by a European Research
496 Council (ERC) Starting Grant to JDJ, as well as Swiss National Science
497 Foundation (FNS) grants to JDJ (31003A_159835) and SV
498 (PMPDP3_158381).

499

500

501 **Conflict of Interest**

502

503 The authors declare no conflict of interest.

504

505 **Data Archiving**

506

507 As a review article, no new data was processed, analyzed, or used directly.

508 **References**

509

510

511 Acevedo A, Brodsky L, Andino R (2014). Mutational and fitness landscapes of
512 an RNA virus revealed through population sequencing. *Nature* **505**: 686-690.

513

514 Anderson RM, May RM (1991). *Infectious diseases of humans: dynamics and*
515 *control*. Oxford University Press: Oxford.

516

517 Bank C, Hietpas RT, Wong A, Bolon DN, Jensen JD (2014). A Bayesian
518 MCMC approach to assess the complete distribution of fitness effects of new
519 mutations: Uncovering the potential for adaptive walks in challenging
520 environments. *Genetics* **196**: 841-852.

521

522 Basdevant A, Goldschmidt C (2008). Asymptotics of the allele frequency
523 spectrum associated with the Bolthausen-Sznitman coalescent. *Electronic*
524 *Journal of Probability* **13**(17): 486-512.

525

526 Begun DJ, Aquadro CF (1992). Levels of naturally occurring DNA
527 polymorphism correlate with recombination rates in *D. melanogaster*. *Nature*
528 **356**: 519-520.

529

530 Berestycki J, Berestycki N, Schweinsberg J (2007). Beta-coalescents and
531 continuous stable random trees. *The Annals of Probability* **35**(5): 1835-1887.

532

533 Berestycki J, Berestycki N, Schweinsberg J (2008). Small-time behavior of
534 beta coalescents. *Annales de l'Institut Henri Poincaré - Probabilités et*
535 *Stastiques* **44**(2): 214-238.

536

537 Bernet GP, Elena SF (2015). Distribution of mutational fitness effects and of
538 epistasis in the 5' untranslated region of a plant RNA virus. *BMC Evolutionary*
539 *Biology* **15**: 274-287.

540

541 Birkner M, Blath J (2008). Computing likelihoods for coalescents with multiple
542 collisions in the infinitely many sites model. *Journal of Mathematical Biology*
543 **57**(3): 435-465.

544

545 Birkner M, Blath J, Eldon B (2013). An ancestral recombination graph for
546 diploid populations with skewed offspring distribution. *Genetics* **193**: 255-290.

547

548 Blath J, Cronjäger MC, Eldon B, Hammer M (2016). The site-frequency
549 spectrum associated with Ξ -coalescents. *Theoretical Population Biology* **110**:
550 36-50.

551

- 552 Bolthausen E, Sznitman AS (1998). On Ruelle's probability cascades and an
553 abstract cavity method. *Communications in Mathematical Physics* **197**: 247-
554 276.
- 555
556 Chare ER, Gould EA, Holmes EC (2003). Phylogenetic analysis reveals a low
557 rate of homologous recombination in negative-sense RNA viruses. *Journal of*
558 *General Virology* **84**: 2691-2703.
- 559
560 Charlesworth B, Morgan MT, Charlesworth D (1993). The effect of deleterious
561 mutations on neutral molecular variation. *Genetics* **134**: 1289-1303.
- 562
563 Comeron JM (2014). Background selection as a baseline for nucleotide
564 variation across the *Drosophila* genome. *PLoS Genetics* **10**(6): e1004434.
- 565
566 Cuevas JM, Geller R, Garijo R, López-Aldeguer J, Sanjuán R (2015).
567 Extremely high mutation rate of HIV-1 in vivo. *PLoS Biology* **13**(9): e1002251.
- 568
569 Davies JL, Simančík F, Lyngsø R, Mailund T, Hein J (2007). On
570 recombination-induced multiple and simultaneous coalescent events.
571 *Genetics* **177**: 2151-2160.
- 572
573 Der R, Epstein CL, Plotkin JB (2011). Generalized population models and the
574 nature of genetic drift. *Theoretical Population Biology* **80**: 80-99.
- 575
576 Donnelly P, Kurtz TG (1999). Particle representations for measure-valued
577 population models. *The Annals of Probability* **27**(1): 166-205.
- 578
579 Eldon B, Birkner M, Blath J, Freund F (2015). Can the site-frequency
580 spectrum distinguish exponential population growth from multiple-merger
581 coalescents? *Genetics* **199**: 841-856.
- 582
583 Eldon B, Degnan JH (2012). Multiple merger gene genealogies in two-
584 species: Monophyly, paraphyly, and polyphyly for two examples of Lambda
585 coalescents. *Theoretical Population Biology* **82**: 117-130.
- 586
587 Eldon B, Wakeley J (2006). Coalescent processes when the distribution of
588 offspring number among individuals is highly skewed. *Genetics* **172**: 2621-
589 2633.
- 590
591 Eldon B, Wakeley J (2008). Linkage disequilibrium under skewed offspring
592 distribution among individuals in a population. *Genetics* **178**: 1517-1532.
- 593
594 Eldon B, Wakeley J (2009). Coalescence times and F_{st} under a skewed
595 offspring distribution among individuals in a population. *Genetics* **181**: 615-
596 629.
- 597
598 Ewens WJ (1979). Testing the generalized neutrality hypothesis. *Theoretical*
599 *Population Biology* **15**(2): 205-216.

600

601 Ewing GB, Jensen JD (2016). The consequences of not accounting for
602 background selection in demographic inference. *Molecular Ecology* **25**: 135-
603 141.

604

605 Eyre-Walker A, Keightley PD (2007). The distribution of fitness effects of new
606 mutations. *Nature Reviews Genetics* **8**: 610-618.

607

608 Ferrer-Admetlla A, Leuenberger C, Jensen JD, Wegmann D (2016). An
609 Approximate Markov Model for the Wright-Fisher Diffusion and its Application
610 to Time Series Data. *Genetics* **203**(2): 831-846.

611

612 Foll M, Poh Y, Renzette N, Ferrer-Admetlla A, Bank C, Shim H *et al* (2014).
613 Influenza virus drug resistance: a time-sampled population genetic
614 perspective. *PLoS Genetics* **10**(2): e1004185.

615

616 Foll M, Shim H, Jensen JD (2015). WFABC: a Wright-Fisher ABC-based
617 approach for inferring effective population sizes and selection coefficients
618 from time-sampled data. *Molecular Ecology Resources* **15**(1): 87-98.

619

620 Grenfell BT, Pybus OG, Gog JR, Wood JLN, Daly JM, Mumford JA *et al*
621 (2004). Unifying the epidemiological and evolutionary dynamics of pathogens.
622 *Science* **303**: 327-332.

623

624 Gutiérrez S, Michalakis Y, Blanc S (2012). Virus population bottlenecks during
625 within-host progression and host-to-host transmission. *Current Opinion in*
626 *Virology* **2**: 546-555.

627

628 Hedgecock D (1994a). Does variance in reproductive success limit effective
629 population sizes of marine organisms? In: Beaumont AR (ed) *Genetics and*
630 *evolution of aquatic organisms*. Chapman & Hall: London, pp 122-133.

631

632 Hedgecock D (1994b). Population genetics of marine organisms. *US Globec*
633 *News* **6**(11): 1-8.

634

635 Hernandez R (2008). A flexible forward simulator for populations subject to
636 selection and demography. *Bioinformatics* **24**(23): 2786-2787.

637

638 Hudson RR, Kaplan NL (1995). Deleterious background selection with
639 recombination. *Genetics* **141**: 1605-1617.

640

641 Huillet T, Möhle M (2011). Population genetics models with skewed fertilities:
642 a forward and backward analysis. *Stochastic Models* **27**: 521-554.

643

644 Innan H, Stephan W (2003). Distinguishing the hitchhiking and background
645 selection models. *Genetics* **165**: 2307-2312.

646

- 647 Jiang L, Liu P, Bank C, Renzette N, Prachanronarong K, Yilmaz LS *et al*
648 (2016). A balance between inhibitor binding and substrate processing confers
649 influenza drug resistance. *Journal of Molecular Biology* **428**: 538-523.
650
- 651 Joseph SB, Swanstrom R (2015). A fitness bottleneck in HIV-1 transmission.
652 *Science* **345**(6193): 136-173.
653
- 654 Kimura M, Maruyama T (1966). The mutational load with epistatic gene
655 interactions in fitness. *Genetics* **54**(6): 1337-1351.
656
- 657 Kingman JFC (1982). The coalescent. *Stochastic Processes and their*
658 *Applications* **13**: 235-248.
659
- 660 Knipe DM, Howley PM (2007). *Fields Virology*, Vol 1. Lippincott Williams &
661 Wilkins: Philadelphia.
662
- 663 Lauring AS, Frydman J, Andino R (2013). The role of mutational robustness in
664 RNA virus evolution. *Nature Reviews Genetics* **11**: 327-336.
665
- 666 Li J, Li H, Jakobsson M, Li S, Sjödin P, Lascoux M (2012). Joint analysis of
667 demography and selection in population genetics: where do we stand and
668 where could we go? *Molecular Ecology* **21**: 28-44.
669
- 670 Li LM, Grassly NC, Fraser C (2014). Genomic analysis of emerging
671 pathogens: methods, application and future trends. *Genome Biology* **15**: 541-
672 550.
673
- 674 Malaspinas A-S (2016). Methods to characterize selective sweeps using time
675 serial samples: an ancient DNA perspective. *Molecular Ecology* **25**: 24-41.
676
- 677 Malaspinas A-S, Malaspinas O, Evans SN, Slatkin M (2012). Estimating allele
678 age and selection coefficient from time-serial data. *Genetics* **192**: 599-607.
679
- 680 McVean GAT, Charlesworth B (2000). The effects of Hill-Robertson
681 interference between weakly selected mutations on patterns of molecular
682 evolution and variation. *Genetics* **155**: 929-944.
683
- 684 Messer PW (2013). SLiM: Simulating evolution with selection and linkage.
685 *Genetics* **194**: 1037-1039.
686
- 687 Möhle M, Sagitov S (2001). A classification of coalescent processes for
688 haploid exchangeable population models. *The Annals of Probability* **29**(4):
689 1547-1562.
690
- 691 Moore CB, John M, James IR, Christiansen FT, Witt CS, Mallal SA (2002).
692 Evidence of HIV-1 adaptation to HLA-restricted immune responses at a
693 population level. *Science* **296**(5572): 1439-1443.
694

- 695 Moran PAP (1958). Random processes in genetics. *Mathematical*
696 *Proceedings of the Cambridge Philosophical Society* **54**(1): 60-71.
697
- 698 Neher RA, Hallatschek O (2013). Genealogies of rapidly adapting populations.
699 *Proceedings of the National Academy of Sciences* **110**(2): 437-442.
700
- 701 Neuhauser C, Krone SM (1997). The genealogy of samples in models with
702 selection. *Genetics* **145**: 519-534.
703
- 704 Nordborg M (1997). Structured coalescent processes on different time scales.
705 *Genetics* **146**: 1501-1514.
706
- 707 Nordborg M, Charlesworth B, Charlesworth D (1996). The effect of
708 recombination on background selection. *Genetical Reserach* **67**(2): 159-174.
709
- 710 Pennings PS, Kryazhimskiy S, Wakeley J (2014). Loss and recovery of
711 genetic diversity in adapting populations of HIV. *PLoS Genetics* **10**(1):
712 e1004000.
713
- 714 Pitman J (1999). Coalescents with multiple collisions. *Journal of Applied*
715 *Probability* **27**: 1870-1902.
716
- 717 Prüfer K, Munch K, Hellmann I, Akagi K, Miller JR, Walenz B *et al* (2012). The
718 bonobo genome compared with the chimpanzee and human genomes. *Nature*
719 **486**: 527-531.
720
- 721 Pybus OG, Rambaut A (2009). Evolutionary analysis of the dynamics of viral
722 infectious disease. *Nature Reviews Genetics* **10**: 540-550.
723
- 724 Renzette N, Gibson L, Bhattacharjee B, Fisher D, Schleiss MR, Jensen JD *et*
725 *al* (2013). Rapid intrahost evolution of human cytomegalovirus is shaped by
726 demography and positive selection. *PLoS Genetics* **9**(9): e1003735.
727
- 728 Renzette N, Kowalik TF, Jensen JD (2016). On the relative roles of
729 background selection and genic hitchhiking in shaping human
730 cytometgalovirus genetic diversity. *Molecular Ecology* **25**(1): 403-413.
731
- 732 Sagitov S (1999). The general coalescent with asynchronous mergers of
733 ancestral lines. *Journal of Applied Probability* **36**: 1116-1125.
734
- 735 Sanjuán R (2010). Mutational fitness effects in RNA and single-stranded DNA
736 viruses: common patterns revealed by site-directed mutagenesis studies.
737 *Philosophical Transactions of the Royal Society B* **365**: 1975-1982.
738
- 739 Schlub TE, Grimm AJ, Smyth RP, Cromer D, Chopra A, Mallal S *et al* (2014).
740 Fifteen to twenty percent of HIV substitution mutations are associated with
741 recombination. *Journal of Virology* **88**(7): 3837-3849.
742

- 743 Schweinsberg J (2000). Coalescents with simultaneous multiple collisions.
744 *Electronic Journal of Probability* **5**(12): 1-50.
745
- 746 Schweinsberg J (2003). Coalescent processes obtained from supercritical
747 Galton-Watson processes. *Stochastic processes and their Applications* **106**:
748 107-139.
749
- 750 Service SK, Blower SM (1995). HIV transmission in sexual networks: an
751 empirical analysis. *Proceedings of the Royal Society of London B: Biological*
752 *Sciences* **260**(1359): 237-244.
753
- 754 Simon-Loriere E, Holmes EC (2011). Why do RNA viruses recombine? *Nature*
755 *Reviews Microbiology* **9**: 617-626.
756
- 757 Steinrücken M, Birkner M, Blath J (2013). Analysis of DNA sequence variation
758 within marine species using Beta-coalescents. *Theoretical Population Biology*
759 **87**: 15-24.
760
- 761 Strelkova N, Lässig M (2012). Clonal interference in the evolution of
762 influenza. *Genetics* **192**: 671-682.
763
- 764 Taylor JE, Véber A (2009). Coalescent processes in subdivided populations
765 subject to recurrent mass extinctions. *Electronic Journal of Probability* **14**(9):
766 242-288.
767
- 768 Tebit DM, Arts EJ (2011). Tracking a century of global expansion and
769 evolution of HIV to drive understanding and to combat disease. *Lancet*
770 *Infectious Disease* **11**: 45-46.
771
- 772 Tellier A, Lemaire C (2014). Coalescence 2.0: a multiple branching of recent
773 theoretical developments and their applications. *Molecular Ecology* **23**: 2637-
774 2652.
775
- 776 Thornton KR (2014). A C++ template library for efficient forward-time
777 population genetic simulation of large populations. *Genetics* **198**: 157-166.
778
- 779 Wilkinson-Herbots HM (1998). Genealogy and subpopulation differentiation
780 under various models of population structure. *Journal of Mathematical Biology*
781 **37**: 535-585.
782
- 783 Williamson EG, Slatkin M (1999). Using maximum likelihood to estimate
784 population size from temporal change in allele frequencies. *Genetics* **152**:
785 755-761.
786
- 787 Wright S (1931). Evolution in Mendelian populations. *Genetics* **16**: 97-159.
788

- 789 Zeng K (2013). A coalescent model of background selection with
790 recombination, demography and variation in selection coefficients. *Heredity*
791 **100**: 363-371.
- 792
793 Zeng K, Charlesworth B (2011). The joint effects of background selection and
794 genetic recombination on local gene genealogies. *Genetics* **189**: 251-266.
- 795
796 Zhu S, Degnan JH, Goldstien SJ, Eldon B (2015). Hybrid-Lambda: simulation
797 of multiple merger and Kingman gene genealogies in species networks and
798 species trees. *BMC Bioinformatics* **16**: 292-298.
- 799
800 Zhu Y, Yongky A, Yin J (2009). Growth of an RNA virus in single cells reveals
801 a broad fitness distribution. *Virology* **385**: 39-46.
- 802
803
804