

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

3

4

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

8

9

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

12 <sup>2</sup> – Swiss Institute of Bioinformatics (SIB), Lausanne, Switzerland

13 <sup>3</sup> – Instituto Gulbenkian de Ciência (IGC), Oeiras, Portugal

14 <sup>4</sup> – Arizona State University (ASU), School of Life Sciences, Center for  
15 Evolution & Medicine, Phoenix, USA

16

17 Word Count: 3400

18

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

21 **Abstract**

22

23 Many features of viral 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 in to 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, and they experience frequent  
42 positive selection in response to host immunity or antiviral treatment.  
43 However, despite these desired attributes, standard population genetic  
44 models must be used with caution when making evolutionary inference.

45

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

57

58 Secondly, the mutation rates of many viruses, particularly RNA viruses, are  
59 among the highest observed across taxa (Lauring *et al*, 2013; Cuevas *et al*,  
60 2015) Though these high rates of mutation are what enable new beneficial  
61 mutations to arise, potentially allowing for rapid resistance to host immunity or  
62 antiviral drugs, they also render high mutational loads (Sanjuán, 2010; Lauring  
63 *et al*, 2013). Specifically, the distribution of fitness effects (DFE) has now been  
64 described across taxa – demonstrating that the input of deleterious mutations  
65 far outnumbers the input of beneficial mutations (Acevedo *et al*, 2014; Bank *et*  
66 *al*, 2014; Bernet and Elena, 2015; Jiang *et al*, 2016). The purging of these  
67 deleterious mutants through purifying selection can affect other areas in the  
68 genome through a process known as background selection (BGS)

69 (Charlesworth *et al*, 1993). Accounting for these effects is important for  
70 accurate evolutionary inference in general (Ewing and Jensen, 2016), but  
71 essential for the study of viruses due to their particularly high rates of mutation  
72 and compact genomes (Renzette *et al*, 2016).

73

74 Given these distinctive features of viral populations and the increasing use of  
75 population genetic inference in this area, it is crucial to account for these  
76 processes that are shaping the amount and distribution of variation across  
77 their genomes. We aim here to draw particular attention to multiple-merger  
78 coalescent events and background selection, and the repercussions of  
79 ignoring them in population genetic inference, highlighting particular  
80 applications to viruses. We conclude with general recommendations for how  
81 best to address these topics in the future.

82

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

84

85 *Inferring evolutionary history using the Wright-Fisher model: benefits and*  
86 *shortcomings*

87

88 Many population genetic statistics and subsequent inference are based on the  
89 Kingman coalescent and the Wright-Fisher (WF) model (Wright, 1931;  
90 Kingman, 1982). With increasing computational power, the WF model has  
91 also been implemented in forward-time methods for inferring population  
92 genetic parameters such as selection coefficients and effective population  
93 sizes ( $N_e$ ) from time-sampled data (Ewens, 1979; Williamson and Slatkin,  
94 1999; Malaspinas *et al*, 2012; Foll *et al*, 2014; Foll *et al*, 2015; Ferrer-Admetlla  
95 *et al*, *in press*). These methods are robust to some violations of WF model  
96 assumptions, such as constant population size, panmixia, and non-  
97 overlapping generations, and also have been extended to accommodate  
98 selection, migration and population structure (Neuhauser and Krone, 1997;  
99 Nordborg, 1997).

100

101 However, it has been suggested that violation of the assumption of a small  
102 variance in offspring number in the WF model leads to erroneous inference of  
103 population genetic parameters (Eldon and Wakeley, 2006). Biological factors  
104 such as sweepstake reproductive events, population bottlenecks, and  
105 recurrent positive selection may lead to skewed distributions in offspring  
106 number (Eldon and Wakeley, 2006; Li *et al*, 2014); examples include various  
107 prokaryotes, fungi, plants, marine organisms, and viruses (reviewed Tellier  
108 and Lemaire, 2014). These skewed offspring distributions can also result in  
109 elevated linkage disequilibrium (LD) despite frequent recombination (Eldon  
110 and Wakeley, 2008; Birkner *et al*, 2013). They may also skew estimates of  $F_{ST}$   
111 from those expected under WF models, as there is a high probability of alleles  
112 being *identical by descent* in subpopulations, where the expectation of  
113 coalescent times within populations is less than that between populations

114 regardless of the timescale or magnitude of gene flow (Eldon and Wakeley,  
115 2009).

116

117

### 118 *Beyond WF assumptions: the Multiple Merger Coalescent*

119

120 A more general coalescent class of models, summarized as the MMC class,  
121 can account for these violations, particularly for (non-Poisson) skewed  
122 offspring distributions, by allowing more than two lineages to coalesce at a  
123 time (Table 1). These are often derived from Moran models, (Moran, 1958),  
124 generalized to allow multiple offspring per individual. In contrast to the WF  
125 model (for which  $P(k > 2) = 0$ ), a probability distribution for  $k$ -merger events  
126 determines coalescence.

127

128 The parameters inferred under the MMC differ notably from those inferred  
129 under the Kingman coalescent in several respects. In a Kingman coalescent,  
130 effective size  $N_e$  scales linearly with census size  $N$ , whereas for the MMC it  
131 does not (Huillet and Möhle, 2011). Thus genetic diversity is only weakly  
132 related to population size. Coalescent trees under the MMC also have more  
133 pronounced star-like genealogies (Figure 1), and their site frequency spectra  
134 (SFSs) are skewed toward an excess of low frequency and high frequency  
135 variants, generating a more negative Tajima's  $D$  (Birkner *et al*, 2013). With  
136 similar migration and population size, alleles fix at a higher rate per population  
137 in the MMC than under the Kingman coalescent, and thus higher  $F_{ST}$  is  
138 expected between subpopulations (Eldon and Wakeley, 2009). Further, the  
139 efficacy of selection increases, as selection acts almost deterministically  
140 between multiple merger events (Tellier and Lemaire, 2014). Theoretical  
141 analyses demonstrate that the fixation probability of a new mutant with a  
142 positive selection coefficient approaches 1 as the population size increases, in  
143 stark contrast with traditional expectations under the standard Wright-Fisher  
144 model (Der *et al*, 2011).

145

146 Not accounting for skewed offspring distributions can lead to mis-inference.  
147 For instance, Eldon and Wakeley (2006) showed that for Pacific oysters,  
148 which have been shown to occasionally undergo sweepstake-like  
149 reproductive events, the estimated population-wide mutation rate  $\theta$  inferred  
150 under the Kingman coalescent is two orders of magnitude larger than that  
151 obtained from the  $\psi$ -coalescent (see below) - 9 vs 0.0308, respectively - and,  
152 indeed, does not well fit the data.

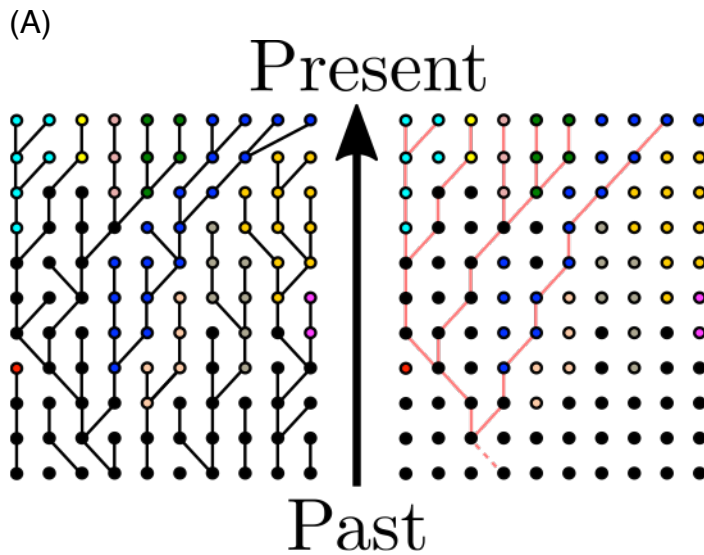
153

154

155 Figure 1: Multiple-Merger and Kingman Coalescent Realizations

156

157

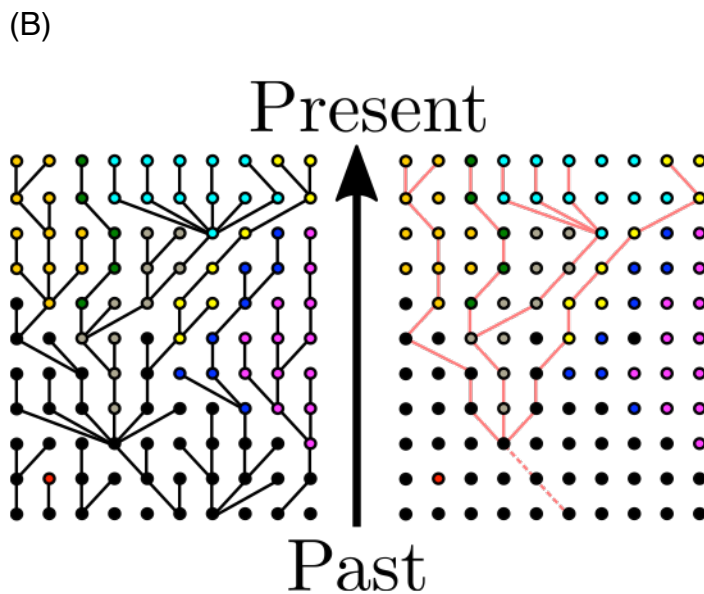


158

159

160

161



162

163

164

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

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

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

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

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

170

171



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

173

Coalescent model	Allows MMs?	Allows simultaneous MMs?	Distribution and parameters	References
$\Xi$ -coalescent	Yes	Yes	MMC events occur with the probability $\Lambda$ , where participating lineages are randomly grouped into $M$ simultaneous mergers with probability $1/M$	Schweinsberg (2000); Möhle and Sagitov (2001)
$\Lambda$ -coalescent	Yes	No	MMC events occur with the probability $\Lambda$ (but $\leq 1$ event/gen)	Donnelly and Kurtz (1999); Pitman (1999); Sagitov (1999)
$\psi$ -coalescent	Yes	No	$\Lambda$ follows a distribution where a fraction $\Psi$ of the population is 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	$\Lambda$ 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	$\Lambda$ 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	$\Lambda$ follows $\beta$ -distribution with $\alpha=2$ : only two lineages are allowed to merge at a time	Kingman (1982)

174

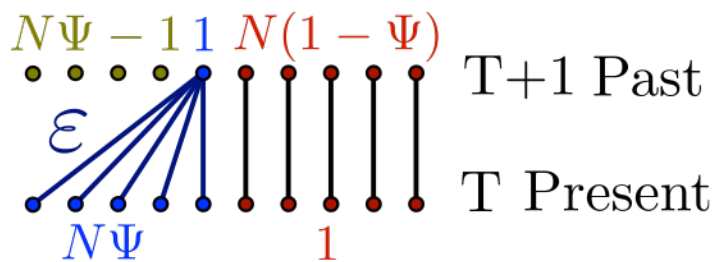
175 *The  $\psi$ -coalescent*

176

177 Introduced by Eldon and Wakeley (2006), the  $\psi$ -coalescent (also called the  
 178 'Dirac-coalescent') differentiates two possible reproductive events (Figure 2)  
 179 (Eldon *et al*, 2015) . Either a standard Moran model reproduction event occurs  
 180 (with probability  $1-\varepsilon$ ), where a single individual is randomly chosen to  
 181 reproduce and the (single) offspring replaces one randomly chosen non-  
 182 parental individual; all other individuals, including the parent, persist.  
 183 Alternatively, a 'sweepstake' reproductive event occurs (with probability  $\varepsilon$ )  
 184 (Hedgcock, 1994)), where a single parent replaces  $\psi*N$  individuals.  
 185 Consequently, an individual may have many offspring and the underlying  
 186 coalescent process will be characterized by MM events, or in the case of the  
 187  $\Xi$ -coalescent, multiple MM events. However, in contrast to other MMC models  
 188 (*e.g.*,  $\lambda$ - or  $\Xi$ -coalescent), the parameter  $\psi$  has a clear biological interpretation  
 189 as the fraction of the population that is replaced in each sweepstake  
 190 reproductive event. Though the assumption of a fixed  $\psi$  (as in the normal  $\psi$ -  
 191 coalescent) seems biologically unrealistic, it can be avoided by treating  $\psi$  as a  
 192 Poisson parameter.

193

194 Figure 2: Depiction of the  $\psi$  coalescent



195

196 Figure 2: Lineages between the present and one generation in the past,  
 197 where  $N$  is the population size,  $\varepsilon$  is the probability of a sweepstake event, and  
 198  $\psi$  is the fraction of the population that is replaced in each such event. Labels  
 199 in the top row give the number of parental individuals reproducing in a given  
 200 manner (represented by color), whereas labels in the bottom row give the  
 201 number of corresponding offspring per parent. Note that time is running  
 202 backwards in the coalescent framework.

203

## 204 *Application to Viruses*

205

206 There are several reasons why a modified Moran model may better capture  
207 viral evolution than models converging to the Kingman coalescent, although it  
208 does not account for fitness differences between individuals. First, viral  
209 evolution is driven by strong bottlenecks during host transmission and  
210 intrahost selection processes, which likely result in skewed offspring  
211 distributions (Figure 3) (Gutiérrez *et al*, 2012; Tellier and Lemaire, 2014).  
212 Further, viruses display the MMC-typical low  $N_e/N$  ratio (Pennings *et al*, 2014;  
213 Tellier and Lemaire, 2014), can adapt rapidly (Neher and Hallatschek, 2013),  
214 and may have sweepstake-like reproductive events in which a single virion  
215 can propagate a large fraction of the entire population (Grenfell *et al*, 2004;  
216 Pybus and Rambaut, 2009). For example, the influenza virus hemagglutinin  
217 (HA) segment appears to be under strong directional selection imposed by  
218 host immunity (and sometimes drug treatment), resulting in a ladder-like  
219 genealogy (Grenfell *et al*, 2004), suggesting that only a few viruses seed the  
220 entire next generation.

221

222 The processes that make viruses ideal candidates for MMCs can, however,  
223 differ by scale (see Figure 3); for example, following transmission events,  
224 there are severe founder events and potentially high recombination within the  
225 host (*e.g.*, HIV, HCMV). Subsequent compartmentalization may introduce  
226 intra-host population structure through bottlenecks, colonization events, and  
227 extinction events (Renzette *et al*, 2013). To date, it remains unclear how often  
228 MMCs fit the patterns of variation observed in intra-host relative to inter-host  
229 viral populations – but such comparisons are increasingly feasible. Finally,  
230 periods of latency - temporary viral inactivation with cessation of reproduction  
231 - should be incorporated in such modeling. Thus, multiple MMC models are a  
232 necessary but not final step towards addressing the various patterns observed  
233 at different scales of virus evolution (Table 1).

234

235 The large data sets often generated from viruses may also prove impractical  
236 for the likelihood-based methods commonly employed for MMCs. This  
237 limitation has partially been overcome by Eldon *et al.* (2015), who proposed  
238 an approximate likelihood method along with an Approximate Bayesian  
239 Computation (ABC) approach based on the SFS to distinguish between the  
240 MMC and exponential population growth. Although both effects are expected  
241 to result in very similar SFSs, characterized by an excess of singletons as  
242 compared to the Kingman coalescent, the bulk and tail of the SFS (i.e., the  
243 higher-order frequency classes) typically differ, which can be assessed by  
244 approximate likelihood-ratio tests and Approximate Bayes Factors (Eldon *et*  
245 *al.*, 2015).

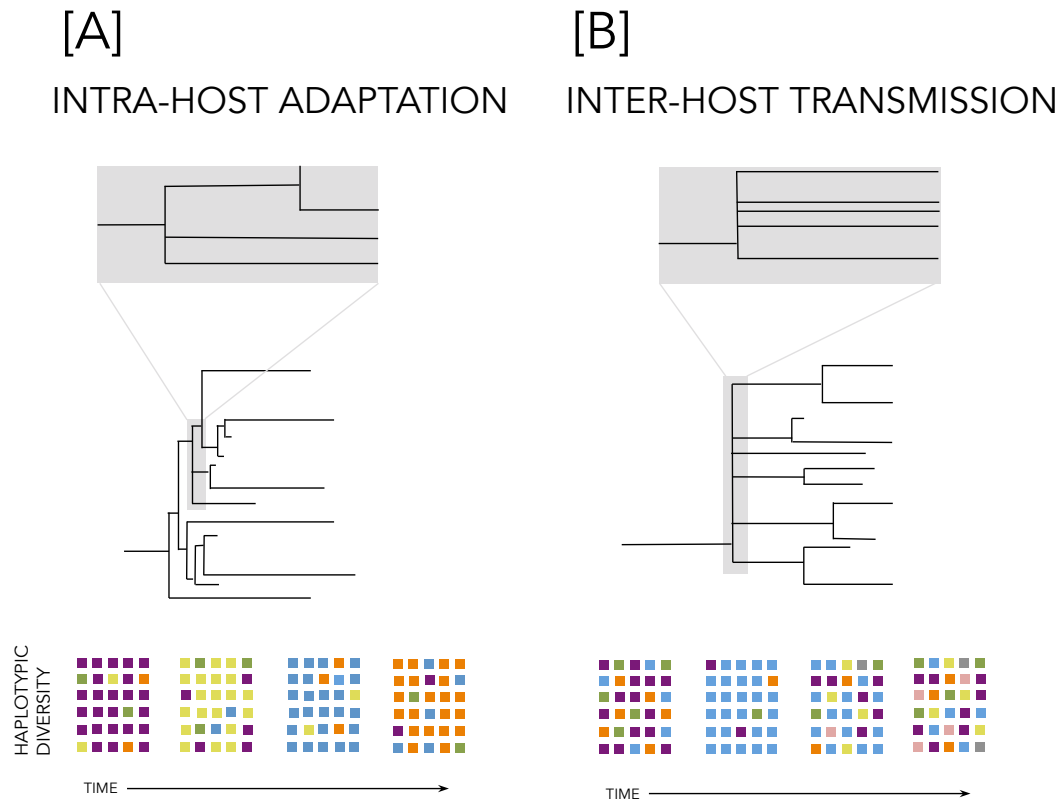
246

247

248

249

### Figure 3: Example Processes Spurring MM Events in Viral Populations



250

251 Figure 3: Examples include (A) intra-host adaptation (a selective process) and  
252 (B) inter-host transmission (a demographic process). The tree in (A)  
253 characterizes, for example, NA or HA evolution in the influenza A virus, driven  
254 by positive selection; selection by host immunity is ongoing, while that from  
255 drug treatment may be intermittent. The tree in (B) represents inter-host  
256 transmission and its associated bottleneck; for viruses that compartmentalize  
257 (such as HCMV and HIV), similar patterns follow transmission to new  
258 compartments. The colored squares below the trees roughly indicate the  
259 diversity of the population through time. Intra-host adaptation may temporarily  
260 decrease diversity owing to genetic hitchhiking, though single snapshots may  
261 not reflect varying temporal levels of diversity. During inter-host transmission,  
262 diversity decreases owing to the associated bottleneck but then may quickly  
263 recover in the new host.

264

265 [ BOX 1: Future challenges in MMC models ]

266

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

269

- 270 1. Describe summary statistics that capture demographic features and  
271 processes when offspring distributions are highly skewed; such  
272 patterns will be required for large-scale inference in a computationally  
273 efficient (e.g., Approximate Bayesian) framework.
- 274 2. Better understand the behavior of commonly used summary statistics  
275 under such models, as done for  $F_{ST}$  by Eldon and Wakeley (2009), for  
276 commonly used divergence, SFS, and LD-based statistics.
- 277 3. Determine which MMCs are best suited for different scales of viral  
278 evolution (i.e., intra-host, inter-host, global); develop novel models if  
279 necessary.
- 280 4. Investigate the effect of violations of MMC assumptions (e.g.,  
281 overlapping generations, number of multiple merger events) on  
282 inference.

283

284 [ END BOX 1 ]

285

## 286 ***Purifying Selection and Linkage in Viral Populations***

287

### 288 *Modeling Background Selection*

289

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

297

298 BGS was originally proposed to explain patterns of reduced diversity in  
299 regions of low recombination – patterns that were previously suggested to be  
300 the signature of genetic hitchhiking (HH) around strongly beneficial mutations  
301 (see Begun and Aquadro 1992 and Charlesworth *et al* 1993). It was argued  
302 that only neutral mutations present on the “least-loaded” chromosomes – that  
303 is, those with the fewest deleterious mutations – have appreciable  
304 probabilities of reaching high frequencies or fixation. Kimura and Maruyama  
305 (1966) showed that the proportion of individuals belonging to the least-loaded  
306 class is

307

$$308 \quad f_0 = \exp\left(-\frac{U}{2hs}\right), \quad (1)$$

309

310 where  $U$  is the rate of mutation to a deleterious state,  $s$  is the selection  
311 coefficient against homozygous mutations, and  $h$  is the dominance coefficient.

312

313 The least-loaded class, and thus genetic diversity in the presence of BGS, is  
314 therefore dependent on the balance between the influx of deleterious  
315 mutations (occurring at rate  $U$ ) and their removal by natural selection  
316 (according to the product  $hs$ ). Assuming that offspring exclusively originate

317 from the least-loaded class of individuals, Charlesworth *et al* (1993)  
318 expressed the expected neutral diversity due to background selection as

319

$$320 \pi = 4 f_o N_e \mu , \quad (2)$$

321

322 where  $N_e$  is the effective population size and  $\mu$  is the mutation rate. As BGS  
323 reduces the number of reproducing individuals, genetic drift increases, thus  
324 reducing genetic diversity and increasing stochasticity in allele trajectories.  
325 Further, since only the genetic diversity segregating in the least-loaded class  
326 can be observed, population size inferred from measures of genetic diversity  
327 may be underestimated if BGS is not properly taken into account (Ewing and  
328 Jensen, 2016).

329

330 In the BGS model described above, strongly deleterious mutations are  
331 maintained in mutation-selection balance such that no skew in the SFS is  
332 expected, as rare variants are rapidly purged. Thus, a simple re-scaling of  $N_e$   
333 is often used as a proxy for the effects of BGS (*e.g.*, Hudson and Kaplan,  
334 1995; Zeng and Charlesworth, 2011; Prüfer *et al*, 2012; Zeng, 2013).

335 However, recent work has demonstrated that, while this re-scaling is  
336 appropriate for strongly deleterious mutations, it is largely inappropriate for  
337 weakly deleterious mutations that may segregate in the population.

338 Experimental work on the shape of the distribution of fitness effects (DFE) in  
339 many organisms indicates that weakly deleterious mutations represent an  
340 important class (*e.g.*, Eyre-Walker and Keightley, 2007; Bank *et al*, 2014).

341 These mutations may act to skew the SFS towards rare alleles as they  
342 decrease the expected frequency of linked neutral mutations relative to  
343 neutral expectations. As subsequent demographic inference is based on the  
344 shape of this SFS, this effect should be properly accounted for by directly  
345 simulating weakly deleterious mutations rather than implementing a simple  
346 rescaling, as is common practice. Figure 4 shows the skew in estimates of  
347 population size and migration rates obtained using an ABC approach when



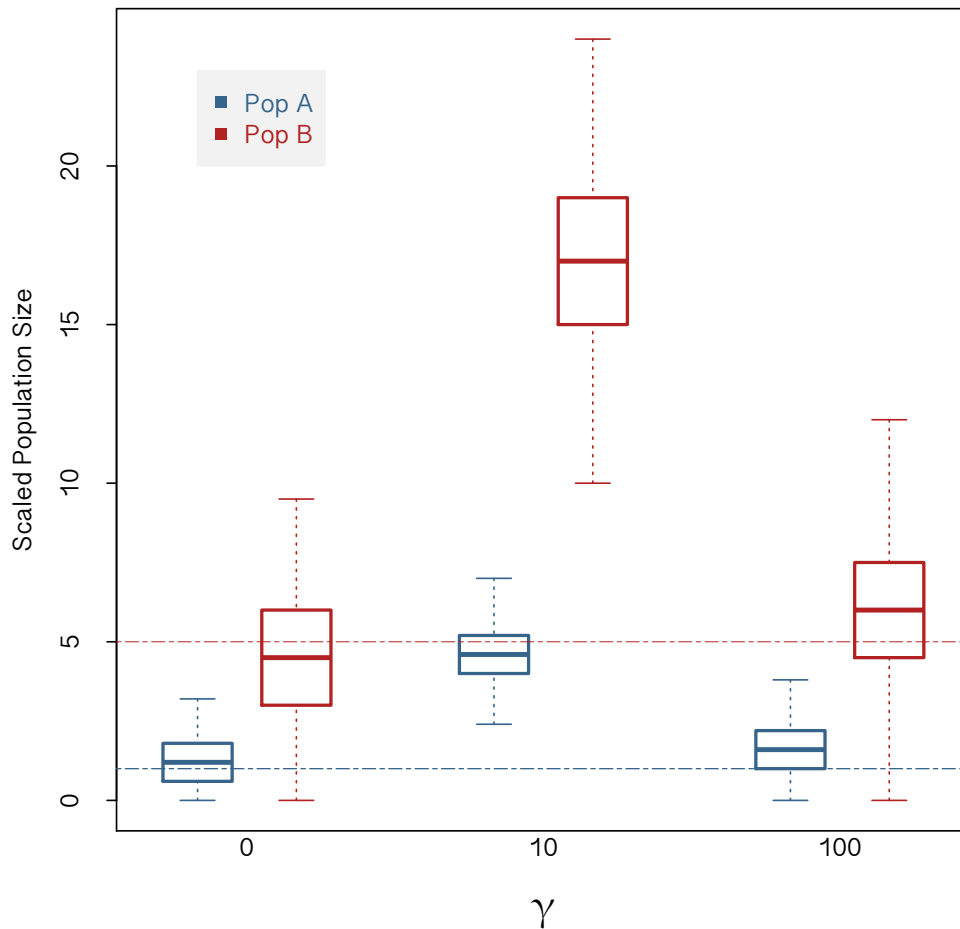
348 BGS is prevalent for two populations A and B that have split at time  $\tau=2N_e$   
349 generations (reproduced from Ewing and Jensen, 2016).

350

351

352

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



354

355 Figure 4: Bias in parameter inference for different levels of BGS, redrawn from  
356 Ewing & Jensen (2016). Posterior densities from ABC inference for population  
357 size are shown. Selection strength is given as  $\gamma$ , where  $\gamma = 2N_e s$ . Population  
358 A has a true scaled size of 1 (blue line), and population B a true scaled size of  
359 5 (red line). As shown, the greatest mis-inference occurs in the presence of  
360 weakly deleterious mutations.

361

362

363

364 *The Effects of Background Selection on Inference for Viral Populations*

365

366 Efforts to estimate the impact of BGS in non-viral organisms have been well  
367 reported. One of the most notable examples is that of Comeron (2014), who  
368 estimated levels of BGS in *Drosophila melanogaster* based on the results of  
369 Hudson and Kaplan (1995) and Nordborg *et al* (1996) using a high-definition  
370 recombination map, with results indicating strong effects across the genome.  
371 For viruses, similar efforts are in their infancy, with the first attempt at such  
372 estimation in a virus reported recently by Renzette *et al* (2016), where they  
373 utilized the theoretical predictions of Innan and Stephan (2003). Interestingly,  
374 the full spectrum of recombination frequencies is available in viral systems –  
375 from non-recombining (*e.g.*, most negative-sense RNA viruses), to re-  
376 assorting (*e.g.*, Influenza virus), to rarely recombining (*e.g.*, Hepatitis C and  
377 West Nile viruses), to frequently recombining (*e.g.*, HIV), offering a highly  
378 promising framework for comparative analyses investigating the  
379 pervasiveness of BGS effects (Chare *et al*, 2003; Simon-Loriere and Holmes,  
380 2011). Further, given the high mutation rates and compact genomes of many  
381 viruses, evolutionary theory suggests effects at least equal to those seen in  
382 *Drosophila*.

383

384 In order to accomplish such inference, improved recombination maps for viral  
385 genomes will be important. With such maps in hand, and given the  
386 amenability of viruses to experimental perturbation, it may indeed be feasible  
387 to understand and account for BGS in models of viral evolution.

388

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

390

391 As BGS almost certainly impacts inference in viral populations, accounting for  
392 its effects is critical. Some future challenges include:

393

- 394 1. Accounting for BGS effects on the SFS by directly simulating weakly  
395 deleterious mutations, rather than by rescaling  $N_e$ .
- 396 2. Improving recombination maps for viral genomes.
- 397 3. Developing models combining the effects of non-equilibrium  
398 demography, positive selection, and BGS, ideally to allow for the joint  
399 estimation of all associated parameters.
- 400 4. Extending methods applied to other taxa to virus populations; for  
401 example, establishing a baseline of variation for use as a null  
402 expectation to estimate BGS levels across the genome, as done for  
403 *Drosophila*.

404

405 [ END BOX 2 ]

406

407

408

409

410

### 411 ***Future Directions***

412

413 Given that skewed offspring distributions and pervasive linked selection are  
414 likely important factors influencing the inference of viral population  
415 parameters, it is important to note that multiple backward and forward  
416 simulation programs have recently been developed which make the modeling  
417 of these processes feasible (Hernandez, 2008; Messer, 2013; Thornton,  
418 2014; Eldon *et al*, 2015; Zhu *et al*, 2015). This will allow researchers to  
419 directly simulate from parameter ranges that may be relevant for their  
420 population of interest, developing a better intuition for the importance of these

421 processes in shaping observed genomic diversity. More concretely, the ability  
422 to now simulate in a computationally efficient framework opens the possibility  
423 of directly implementing ABC inference approaches under these models.  
424 Thus, by drawing mutations from a biologically realistic distribution of fitness  
425 effects and allowing offspring distributions to appropriately vary, it is now  
426 possible to re-implement common demographic estimation or genome scan  
427 approaches; such modified approaches would be based on more appropriate  
428 null expectations of the shape of the SFS, the extent of linkage disequilibrium,  
429 and the degree of population divergence.

430 **Acknowledgements**

431

432 We would like to thank Bjarki Eldon for helpful suggestions during the early  
433 stages of this manuscript. This work was funded by a European Research  
434 Council (ERC) Starting Grant to JDJ, as well as Swiss National Science  
435 Foundation (FNS) grants to JDJ (31003A\_159835) and SV  
436 (PMPDP3\_158381).

437

438

439 **Conflict of Interest**

440

441 The authors declare no conflict of interest.

442

443 **Data Archiving**

444

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

446 **References**

447

448

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

451

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

456

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

460

461 Begun DJ, Aquadro CF (1992). Levels of naturally occurring DNA  
462 polymorphism correlate with recombination rates in *D. melanogaster*. *Nature*  
463 **356**: 519-520.

464

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

467

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

471

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

475

476 Birkner M, Blath J (2008). Computing likelihoods for coalescents with multiple  
477 collisions in the infinitely many sites model. *Journal of Mathematical Biology*  
478 **57**(3): 435-465.

479

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

482

483 Bolthausen E, Sznitman AS (1998). On Ruelle's probability cascades and an  
484 abstract cavity method. *Communications in Mathematical Physics* **197**: 247-  
485 276.

486

487 Chare ER, Gould EA, Holmes EC (2003). Phylogenetic analysis reveals a low  
488 rate of homologous recombination in negative-sense RNA viruses. *Journal of  
489 General Virology* **84**: 2691-2703.

490

- 491 Charlesworth B, Morgan MT, Charlesworth D (1993). The effect of deleterious  
492 mutations on neutral molecular variation. *Genetics* **134**: 1289-1303.
- 493
- 494 Comeron JM (2014). Background selection as a baseline for nucleotide  
495 variation across the *Drosophila* genome. *PLoS Genetics* **10**(6): e1004434.
- 496
- 497 Cuevas JM, Geller R, Garijo R, López-Aldeguer J, Sanjuán R (2015).  
498 Extremely high mutation rate of HIV-1 in vivo. *PLoS Biology* **13**(9): e1002251.
- 499
- 500 Der R, Epstein CL, Plotkin JB (2011). Generalized population models and the  
501 nature of genetic drift. *Theoretical Population Biology* **80**: 80-99.
- 502
- 503 Donnelly P, Kurtz TG (1999). Particle representations for measure-valued  
504 population models. *The Annals of Probability* **27**(1): 166-205.
- 505
- 506 Eldon B, Birkner M, Blath J, Freund F (2015). Can the site-frequency  
507 spectrum distinguish exponential population growth from multiple-merger  
508 coalescents? *Genetics* **199**: 841-856.
- 509
- 510 Eldon B, Degnan JH (2012). Multiple merger gene genealogies in two-  
511 species: Monophyly, paraphyly, and polyphyly for two examples of Lambda  
512 coalescents. *Theoretical Population Biology* **82**: 117-130.
- 513
- 514 Eldon B, Wakeley J (2006). Coalescent processes when the distribution of  
515 offspring number among individuals is highly skewed. *Genetics* **172**: 2621-  
516 2633.
- 517
- 518 Eldon B, Wakeley J (2008). Linkage disequilibrium under skewed offspring  
519 distribution among individuals in a population. *Genetics* **178**: 1517-1532.
- 520
- 521 Eldon B, Wakeley J (2009). Coalescence times and  $F_{st}$  under a skewed  
522 offspring distribution among individuals in a population. *Genetics* **181**: 615-  
523 629.
- 524
- 525 Ewens WJ (1979). Testing the generalized neutrality hypothesis. *Theoretical*  
526 *Population Biology* **15**(2): 205-216.
- 527
- 528 Ewing GB, Jensen JD (2016). The consequences of not accounting for  
529 background selection in demographic inference. *Molecular Ecology* **25**: 135-  
530 141.
- 531
- 532 Eyre-Walker A, Keightley PD (2007). The distribution of fitness effects of new  
533 mutations. *Nature Reviews Genetics* **8**: 610-618.
- 534
- 535 Ferrer-Admetlla A, Leuenberger C, Jensen JD, Wegmann D (*in press*). An  
536 Approximate Markov Model for the Wright-Fisher Diffusion and its Application  
537 to Time Series Data. *Genetics*.
- 538

- 539 Foll M, Poh Y, Renzette N, Ferrer-Admetlla A, Bank C, Shim H *et al* (2014).  
540 Influenza virus drug resistance: a time-sampled population genetic  
541 perspective. *PLoS Genetics* **10**(2): e1004185.  
542
- 543 Foll M, Shim H, Jensen JD (2015). WFABC: a Wright-Fisher ABC-based  
544 approach for inferring effective population sizes and selection coefficients  
545 from time-sampled data. *Molecular Ecology Resources* **15**(1): 87-98.  
546
- 547 Grenfell BT, Pybus OG, Gog JR, Wood JLN, Daly JM, Mumford JA *et al*  
548 (2004). Unifying the epidemiological and evolutionary dynamics of pathogens.  
549 *Science* **303**: 327-332.  
550
- 551 Gutiérrez S, Michalakis Y, Blanc S (2012). Virus population bottlenecks during  
552 within-host progression and host-to-host transmission. *Current Opinion in*  
553 *Virology* **2**: 546-555.  
554
- 555 Hedgecock D (1994). Population genetics of marine organisms. *US Globec*  
556 *News* **6**(11): 1-8.  
557
- 558 Hernandez R (2008). A flexible forward simulator for populations subject to  
559 selection and demography. *Bioinformatics* **24**(23): 2786-2787.  
560
- 561 Hudson RR, Kaplan NL (1995). Deleterious background selection with  
562 recombination. *Genetics* **141**: 1605-1617.  
563
- 564 Huillet T, Möhle M (2011). Population genetics models with skewed fertilities:  
565 a forward and backward analysis. *Stochastic Models* **27**: 521-554.  
566
- 567 Innan H, Stephan W (2003). Distinguishing the hitchhiking and background  
568 selection models. *Genetics* **165**: 2307-2312.  
569
- 570 Jiang L, Liu P, Bank C, Renzette N, Prachanronarong K, Yilmaz LS *et al*  
571 (2016). A balance between inhibitor binding and substrate processing confers  
572 influenza drug resistance. *Journal of Molecular Biology* **428**: 538-523.  
573
- 574 Kimura M, Maruyama T (1966). The mutational load with epistatic gene  
575 interactions in fitness. *Genetics* **54**(6): 1337-1351.  
576
- 577 Kingman JFC (1982). The coalescent. *Stochastic Processes and their*  
578 *Applications* **13**: 235-248.  
579
- 580 Knipe DM, Howley PM (2007). *Fields Virology*, Vol 1. Lippincott Williams &  
581 Wilkins: Philadelphia.  
582
- 583 Lauring AS, Frydman J, Andino R (2013). The role of mutational robustness in  
584 RNA virus evolution. *Nature Reviews Genetics* **11**: 327-336.  
585



- 586 Li J, Li H, Jakobsson M, Li S, Sjödin P, Lascoux M (2012). Joint analysis of  
587 demography and selection in population genetics: where do we stand and  
588 where could we go? *Molecular Ecology* **21**: 28-44.
- 589  
590 Li LM, Grassly NC, Fraser C (2014). Genomic analysis of emerging  
591 pathogens: methods, application and future trends. *Genome Biology* **15**: 541-  
592 550.
- 593  
594 Malaspinas A, Malaspinas O, Evans SN, Slatkin M (2012). Estimating allele  
595 age and selection coefficient from time-serial data. *Genetics* **192**: 599-607.
- 596  
597 Messer PW (2013). SLiM: Simulating evolution with selection and linkage.  
598 *Genetics* **194**: 1037-1039.
- 599  
600 Möhle M, Sagitov S (2001). A classification of coalescent processes for  
601 haploid exchangeable population models. *The Annals of Probability* **29**(4):  
602 1547-1562.
- 603  
604 Moran PAP (1958). Random processes in genetics. *Mathematical*  
605 *Proceedings of the Cambridge Philosophical Society* **54**(1): 60-71.
- 606  
607 Neher RA, Hallatschek O (2013). Genealogies of rapidly adapting populations.  
608 *Proceedings of the National Academy of Sciences* **110**(2): 437-442.
- 609  
610 Neuhauser C, Krone SM (1997). The genealogy of samples in models with  
611 selection. *Genetics* **145**: 519-534.
- 612  
613 Nordborg M (1997). Structured coalescent processes on different time scales.  
614 *Genetics* **146**: 1501-1514.
- 615  
616 Nordborg M, Charlesworth B, Charlesworth D (1996). The effect of  
617 recombination on background selection. *Genetical Research* **67**(2): 159-174.
- 618  
619 Pennings PS, Kryazhimskiy S, Wakeley J (2014). Loss and recovery of  
620 genetic diversity in adapting populations of HIV. *PLoS Genetics* **10**(1):  
621 e1004000.
- 622  
623 Pitman J (1999). Coalescents with multiple collisions. *Journal of Applied*  
624 *Probability* **27**: 1870-1902.
- 625  
626 Prüfer K, Munch K, Hellmann I, Akagi K, Miller JR, Walenz B *et al* (2012). The  
627 bonobo genome compared with the chimpanzee and human genomes. *Nature*  
628 **486**: 527-531.
- 629  
630 Pybus OG, Rambaut A (2009). Evolutionary analysis of the dynamics of viral  
631 infectious disease. *Nature Reviews Genetics* **10**: 540-550.
- 632

- 633 Renzette N, Gibson L, Bhattacharjee B, Fisher D, Schleiss MR, Jensen JD *et*  
634 *al* (2013). Rapid intrahost evolution of human cytomegalovirus is shaped by  
635 demography and positive selection. *PLoS Genetics* **9**(9): e1003735.  
636
- 637 Renzette N, Kowalik TF, Jensen JD (2016). On the relative rules of  
638 background selection and geneic hitchhiking in shaping human  
639 cytometgalovirus genetic diversity. *Molecular Ecology* **25**(1): 403-413.  
640
- 641 Sagitov S (1999). The general coalescent with asynchronous mergers of  
642 ancestral lines. *Journal of Applied Probability* **36**: 1116-1125.  
643
- 644 Sanjuán R (2010). Mutational fitness effects in RNA and single-stranded DNA  
645 viruses: common patterns revealed by site-directed mutagenesis studies.  
646 *Philosophical Transactions of the Royal Society B* **365**: 1975-1982.  
647
- 648 Schweinsberg J (2000). Coalescents with simultaneous multiple collisions.  
649 *Electronic Journal of Probability* **5**(12): 1-50.  
650
- 651 Schweinsberg J (2003). Coalescent processes obtained from supercritical  
652 Galton-Watson processes. *Stochastic processes and their Applications* **106**:  
653 107-139.  
654
- 655 Simon-Loriere E, Holmes EC (2011). Why do RNA viruses recombine? *Nature*  
656 *Reviews Microbiology* **9**: 617-626.  
657
- 658 Steinrücken M, Birkner M, Blath J (2013). Analysis of DNA sequence variation  
659 within marine species using Beta-coalescents. *Theoretical Population Biology*  
660 **87**: 15-24.  
661
- 662 Tellier A, Lemaire C (2014). Coalescence 2.0: a multiple branching of recent  
663 theoretical developments and their applications. *Molecular Ecology* **23**: 2637-  
664 2652.  
665
- 666 Thornton KR (2014). A C++ template library for efficient forward-time  
667 population genetic simulation of large populations. *Genetics* **198**: 157-166.  
668
- 669 Williamson EG, Slatkin M (1999). Using maximum likelihood to estimate  
670 population size from temporal change in allele frequencies. *Genetics* **152**:  
671 755-761.  
672
- 673 Wright S (1931). Evolution in Mendelian populations. *Genetics* **16**: 97-159.  
674
- 675 Zeng K (2013). A coalescent model of background selection with  
676 recombination, demography and variation in selection coefficients. *Heredity*  
677 **100**: 363-371.  
678
- 679 Zeng K, Charlesworth B (2011). The joint effects of background selection and  
680 genetic recombination on local gene genealogies. *Genetics* **189**: 251-266.

681

682 Zhu S, Degnan JH, Goldstien SJ, Eldon B (2015). Hybrid-Lambda: simulation  
683 of multiple merger and Kingman gene genealogies in species networks and  
684 species trees. *BMC Bioinformatics* **16**: 292-298.

685

686

687

688

689