**The effective population size and mutation rate of influenza A virus in acutely infected individuals**

Running Title: Within-host model of influenza

John T. McCrone [1†], Robert J. Woods [2], Arnold S. Monto [3], Emily T. Martin [3], and Adam S. Lauring [1,2 *]

[1] Department of Microbiology and Immunology, University of Michigan, Ann Arbor, MI 48109

[2] Division of Infectious Diseases, Department of Internal Medicine, University of Michigan, Ann Arbor, MI 48109

[3] Department of Epidemiology, University of Michigan, Ann Arbor, MI 48109

[†] Present address

Institute of Evolutionary Biology, University of Edinburgh, Edinburgh EH9; 3FL, UK.

* Corresponding author

Adam S. Lauring

1150 W. Medical Center Dr.

MSRB1 Room 5510B

Ann Arbor, MI 48109-5680

(734) 764-7731

alauring@med.umich.edu

Key Words – influenza virus, transmission, diversity, evolution

1    **Abstract**

2    The global evolutionary dynamics of influenza viruses ultimately derive from processes that take

3    place within and between infected individuals. Recent work suggests that within-host

4    populations are dynamic, but an *in vivo* estimate of mutation rate and population size in

5    naturally infected individuals remains elusive. Here we model the within-host dynamics of

6    influenza A viruses using high depth of coverage sequence data from 200 acute infections in an

7    outpatient, community setting. Using a Wright-Fisher model, we estimate a within-host effective

8    population size of 32-72 and an *in vivo* mutation rate of $3.4 \times 10^{-6}$ per nucleotide per generation.

9

10   **Introduction**

11   The rapid evolution of influenza viruses places demographic processes such as population

12   growth, transmission, and epidemiological spread on a similar time scale as the accumulation of

13   genetic substitutions. This similarity of scale makes it possible to infer demographic processes

14   from genetic sequence data using phylodynamic methods (Lemey et al. 2009; Bedford et al.

15   2014; Bedford et al. 2015). Investigations of the global dynamics of influenza have been

16   successful, in part, because the complexities of within- and between-host processes can be

17   collapsed into a limited number of parameters in the coalescent or birth-death process when

18   averaged over large spatial and temporal scales. However, it becomes increasingly important to

19   disentangle these processes to address more granular questions; for example, the transmission

20   of viruses at local scales or selective pressures imposed by vaccines, antivirals, or novel hosts.

21

22   Phylogenetic approaches that separate within-host processes from those acting at

23   epidemiological scales rely on simple population genetic models to capture the complex

24   dynamics that occur within infected individuals (Didelot et al. 2014; Hall et al. 2015; Didelot et al.

25   2017; De Maio et al. 2018). However, the accuracy of these models depends on reliable

26   estimates of the within-host effective population size ($N_e$), which in the case of influenza virus,

27    has proven difficult due to inherent challenges in collecting longitudinal samples from

28    representative infections. Here, we take advantage of a large, well-studied, community cohort

29    with robust deep sequencing data, from which two important results have emerged (McCrone et

30    al. 2018). First, within-host selection for novel antigenic variants is weak, and second,

31    transmission between hosts imposes a significant bottleneck on the viral population. We

32    leverage these findings to fit a Wright-Fisher model to capture the dynamics of within-host

33    populations. This model provides consistent and robust estimates of the within-host $N_e$ and

34    mutation rate of influenza A virus (IAV) when applied to cross-sectional and longitudinal

35    samples. These findings provide an important baseline for defining processes related to the

36    local dynamics of IAV, and of RNA viruses in general.

37

38    **Results**

39    We recently performed high depth of coverage sequencing of 249 IAV populations recovered

40    from 200 individuals enrolled in the Household Influenza Vaccine Effectiveness (HIVE) study

41    (McCrone et al. 2018). This large number of samples collected within a prospective community-

42    based cohort is a rich dataset for exploring influenza virus evolution over the course of a natural

43    infection. In this and other works, we have documented our sensitivity and specificity for

44    detection of intrahost single nucleotide variants (iSNV) and our error in allele frequency

45    measurement (McCrone and Lauring 2016; Debbink et al. 2017; McCrone et al. 2018). Our

46    dataset also includes 49 serially sampled individuals, who provided a self-collected specimen at

47    the time of symptom onset and a clinic-collected specimen 0–7 days later. This affords an

48    opportunity to explore changes in viral populations in naturally infected individuals over a short

49    time scale.

50

51    We applied a continuous diffusion approximation of the Wright-Fisher model to define the within-

52    host accumulation of mutations using 196 cross-sectional samples, collected 1-7 days following

53    the onset of symptoms (Rouzine et al. 2001). Because we have previously estimated an

54    effective transmission bottleneck of 1-2 genetically distinct variants, we made the simplifying

55    assumption that each infection was clonal and modeled the accumulation of diversity until the

56    time of sampling as a neutral process. Maximum likelihood optimization of this model estimated

57    an *in vivo* mutation rate of $3.4 \times 10^{-6}$ (95% CI $3.1\text{-}3.7 \times 10^{-6}$) mutations per nucleotide per

58    generation (6 hours) and a within-host $N_e$ of 36 ( 95% CI 31-41, Figure 1). We have recently

59    estimated that the majority of mutations in IAV are detrimental and therefore unlikely to be

60    observed at detectable frequencies (Visher et al. 2016). As only ~10% of mutations in influenza

61    A virus are neutral, we propose that the true *in vivo* mutation rate is approximately ten-fold

62    higher than our estimated rate, which does not account for purifying selection. This results in an

63    *in vivo* mutation rate of approximately $3.4 \times 10^{-5}$ substitutions per nucleotide replicated per

64    generation, which is within the range of estimates for IAV's biochemical mutation rate in

65    epithelial cells (Sanjuán et al. 2010).

66

67    To determine the robustness of our $N_e$ estimate, we fit this same model to changes in allele

68    frequencies observed in a subset of paired longitudinal samples. We restricted this analysis to

69    alleles observed at the first time point in samples taken at least 1 day apart (63 iSNV in 29

70    sample pairs). There was very little change in iSNV frequency in populations sampled twice on

71    the same day ($R^2 = 0.986$, Figure 2, Supplement 1A of (McCrone et al. 2018)). The

72    concordance of same-day samples suggests that our sampling procedure and frequency

73    measurements are reproducible. Maximum likelihood optimization of this model revealed a

74    within-host $N_e$ of 34 (95% CI 25-46, Table 1), very similar to that observed in the cross-sectional

75    data above. Comparable estimates were obtained when synonymous and nonsynomous

76    mutations were fit separately (Table 1). As there is some uncertainty in the within-host

77    generation time (Geoghegan et al. 2016), we also estimated the $N_e$ based on a 12 hour

78    generation. As expected, increasing the generation time results in a smaller $N_e$.

3

79

80    The Wright-Fisher model assumes that each allele in a population is independent. This

81    assumption would be violated if there were multiple iSNV per genomic segment or varying

82    linkage of iSNV across segments due to reassortment. However, heterotypic reassortment is

83    quite rare within hosts (Sobel Leonard et al. 2017), and the per-sample diversity in our dataset

84    was sufficiently low that nearly all segments had either 0 or 1 iSNV. To ensure that our results

85    were robust to the assumption of independent allele frequencies, we fit the above model 500

86    times, each time randomly subsetting our data such that only one iSNV per individual was

87    included. In practice, this approach also tested the sensitivity of our estimates to individual allele

88    trajectories. Under these conditions, we found a median $N_e$ of 42 (IQR 37-52, Figure 1B). Thus,

89    in the initial analysis, non-independence among iSNV within the same host may have caused a

90    slight bias due to a few hosts with extreme frequency changes.

91

92    The estimates above include the probability that undetected variants are present but missed

93    due to imperfect sensitivity (see Methods and (McCrone and Lauring 2016)); however, they do

94    not account for uncertainty in the frequency measurements, which if large, would bias the $N_e$

95    estimate toward lower values. To accommodate this uncertainty we relied on the fact that 141 of

96    the 249 samples in were amplified and sequenced in duplicate (McCrone et al. 2018). We

97    modeled the frequency-dependent variance present in the data as a beta distribution with $\alpha =$

98    $p * n, \ \beta = p(1 - p) * n$, where $p$ represents the true frequency (the mean in the duplicate

99    measurements) and $n$ roughly represents the number of samples in a binomial distribution with

100    probability $p$, and was determined with maximum likelihood optimization. We then adapted a

101    Bayesian approach and estimated the posterior distribution of $N_e$ integrated over all

102    unobserved, true frequency trajectories. The analysis resulted in a marginally increased $N_e$

103    estimate of 50 (32-72 95% HPD, Figure 1C). The agreement between this model and our

104    previous estimates suggests that the relatively small $N_e$ is driven by the allele trajectories

105    themselves and is not the result of uncertainty in our frequency measurements.

106

107    **Discussion**

108    We have investigated the within-host dynamics of influenza in a large, well-defined cohort of

109    representative infections and found that, under a Wright-Fisher model, the population is

110    characterized by a small effective population size. Our findings differ from those reported in

111    studies of immunosuppressed, chronically infected individuals, which have shown that within-

112    host populations of influenza virus are characterized by large effective population sizes, clonal

113    interference, and selective pressures that mimic those seen at larger biological scales (Xue et

114    al. 2017; Lumby et al. 2020). The difference in these $N_e$ estimates likely lies in the fundamental

115    difference between the population dynamics of acute and chronic infections. Chronic infections,

116    which manifest in rare immunologically atypical hosts, establish large, stable populations and

117    may be "insulated" from the drastic fluctuation in population size that define acute cases. In the

118    absence of any evidence for antigenic selection, it seems that evolution during the early period

119    of influenza infections, the time frame during which transmission is most likely to occur, is best

120    modelled as a stochastic process.

121

122    The Wright-Fisher model provides a simple framework for exploring the evolutionary dynamics

123    of "real-world" populations. The model's tractability comes at the cost of many simplifying

124    assumptions (e.g. constant population size, discrete generations, homogenous mixing, neutral

125    evolution), which are rarely, if ever, met by biological populations. Influenza viruses clearly exist

126    as complex populations whose evolutionary dynamics are influenced by a mixture of processes

127    not captured explicitly in the Wright-Fisher model (e.g. deleterious mutation load, migration

128    between sites of infection, rapid population growth and decline (Lakdawala et al. 2015; Visher et

129    al. 2016; Zhao et al. 2019)). However, the detailed, longitudinal sampling needed to fit models

130    that explicitly capture this complexity is not available for most influenza infections, which are

131    typically short-lived and not medically attended. In the absence of such data, we have chosen a

132    more tractable model that can yield reliable estimates of the general tendencies, rather than

133    more complex models that may lack identifiability and generalizability.

134

135    These estimates of the effective population size and mutation rate, combined with previous

136    estimates of the transmission bottleneck, provide a useful expectation for the shared diversity

137    between direct transmission pairs, and can be used in conjunction with standard

138    epidemiological models to study the forces that drive influenza evolution at a granular level.

139

140    **Methods**

141

142    *Fitting mutation rate and $N_e$*

143    The diffusion approximation to the Wright-Fisher model makes predictions regarding the allele

144    frequency spectrum of a population given a mutation rate and $N_e$. Starting from a monomorphic

145    state, while t<<$N_e$, the probability of observing a mutation at frequency $p_t$ be approximated as in

146    equation 85 of (Rouzine et al. 2001)

147

148    $P(p_t, | \, t, \mu, N_e) = \frac{2\mu N_e}{p_t} e^{-\frac{2N_e p_t}{t}}$                                                                                      (1)

149    Where $\mu$ is the mutation rate in substitutions/site/generation, $N_e$ is the effective population size

150    and $t$ is the number of generations. Consistent with previous models of within-host influenza, we

151    set the generation time to 6 hours (Geoghegan et al. 2016). We further assumed that infection

152    began 1 day prior to symptom onset (Carrat et al. 2008).

153

6

154    To account for limitations in iSNV detection, we integrated over regions of the probability density

155    where we have observed less than perfect sensitivity. The probability of not observing an iSNV

156    at a locus is given by summing over the possibilities that (i) a mutation is present but below our

157    level of detection $P(p_t \approx 0 \mid p_t < 0.02, t, \mu, N_e)$, and (ii) a mutation is present but missed

158    due to low sensitivity at low frequencies $P(p_t \approx 0 \mid 0.02 < p_t < 0.1, t, \mu, N_e)$. In this model,

159    we assumed there were 13,133 polymorphic loci in each sample (the number of coding sites

160    present in the reference strain from 2014-2015). Under these assumptions,

161

$$
\begin{aligned}
162 \quad & P(p_t \approx 0 \mid t, \mu, N_e) = \\
163 \quad & P(p_t \approx 0 \mid p_t < 0.02, t, \mu, N_e) + \\
164 \quad & P(p_t \approx 0, t \mid 0.02 < p_t < 0.1, t, \mu, N_e)
\end{aligned}
$$

165    (2)

166    Where

167    $P(p_t \approx 0 \mid p_t < 0.02, t, \mu, N_e) = \int_0^{0.02} P(p_t, \mid t, \mu, N_e) dp_t$     (3)

168    and

169    $P(p_t \approx 0 \mid 0.02 < p_t < 0.1, t, \mu, N_e) = \sum_{f_i}^{[0.02, 0.05, 0.10)} (FNR \mid Titer_r, f_i) \int_{f_i}^{f_{i+1}} P(p_t \mid$

170    $\mu, t, N_e) dp_t$     (4)

171

172    Where $(FNR \mid Titer_r, f_i)$ is the false negative rate given the frequency and the sample titer

173    (See Supplementary File 1 in (McCrone et al. 2018)). As before, we assumed the sensitivity in

174    the intervals between 0.02, 0.05 and 0.1 was equal to the sensitivity at the lower bound, and

175    that the sensitivity was perfect at frequencies above 0.1. The log-likelihood of a given $\mu$ and $N_e$

176    pair is then the sum of the log of equations 1 and 2 for all possible sites in the data set. The

177    maximum-likelihood values were estimated using the bbmle package in R (Ben Bolker and

178    Team 2020; Team 2020).

179

180 *Diffusion approximation*

181 We implemented the diffusion approximation as in (Kimura 1955), with minor modifications. As

182 above, we included the limitations in our sensitivity to detect rare iSNV by integrating over all

183 possible explanations for why an iSNV might not be observed at the second time point.

184

185 *Bayesian implementation of the diffusion approximation*

186 To account for measurement error in our estimates we adopted a similar approach to that

187 developed in (Williamson and Slatkin 1999). The likelihood of observing frequencies $\widehat{p_0}, \widehat{p_t}$ at

188 time 0 and t given the true frequencies $p_0$ and $p_t$

$$P\left(\widehat{p_0}, \widehat{p_t} | N_e, p_0, p_t\right) = P(\widehat{p_0}|p_0)P(p_t|p_0, N_e)P(\widehat{p_t}|p_t)$$

190                                                                                                             (5)

191 where $P(\widehat{p_x}|p_x)$ accounts for measurement error and is defined for $\widehat{p_x} > 0$ by the probability

192 density at $\widehat{p_x}$ of a beta distribution with $\alpha = p * n$, $\beta = p(1-p) * n$ where n=503 and was

193 determined from the estimating the error in replicate sequencing samples.

194 In cases where $\widehat{p_x} = 0$ and $p_x > 0$, $P(\widehat{p_x}|p_x)$ is the sum of the cumulative density function of the

195 same beta distribution up to 0.02 (i.e. the variant is detected below the limit of detection) and the

196 probability of not detecting the variant given the sample titer and false negative rate as above

197 (the variant was not observed to imperfect sensitivity). $P(p_t|p_0, N_e)$ is the transition probability of

198 a variant at frequency $p_0$ to drifting to $p_t$ given t generations and an the effective population size

199 of $N_e$ as in equation 15' in (Kimura 1955). The posterior is proportional to the product of this

200 likelihood and priors on $N_e$, $p_0$, and $p_t$. We choose uniform priors for $p_0$ and $p_t$ and a diffuse

201 gamma prior with shape of 0.036 and scale of 1000 (mean 36 as informed by the cross-

202 sectional data analysis). As with the other analyses the generation time was set to 6 hours

203 (Geoghegan et al. 2016). This approach was implemented as a plugin for BEAST and the

204 posterior was estimated using BEAST v1.10.4 (Suchard et al. 2018). Ten independent MCMC

8

205    chains were run for 10 million states. Each chain was sampled every 10,000 iterations with the

206    first 1 million states discarded as burn in. All ten chains were combined and ESS for all

207    parameters was >200.  Convergence was assessed in *TRACER* (Rambaut et al. 2018).

208

209    **Acknowledgments**

215

216    **Data availability**

217    All raw sequence data have been deposited at the NCBI sequence read archive (BioProject

218    Accession number PRJNA412631) as described in (McCrone et al. 2018). Variants were called

219    following the validated protocol outlined in (McCrone and Lauring 2016) with details provided in

220    (McCrone et al. 2018). Called variants and the scripts needed to reproduce this analysis are

221    publicly available at https://github.com/lauringlab/IAV_within-host_Ne

222

223    **References**

224    Bedford T, Riley S, Barr IG, Broor S, Chadha M, Cox NJ, Daniels RS, Gunasekaran CP, Hurt
225        AC, Kelso A, et al. 2015. Global circulation patterns of seasonal influenza viruses vary with
226        antigenic drift. Nature 523:217–220.

227    Bedford T, Suchard MA, Lemey P, Dudas G, Gregory V, Hay AJ, McCauley JW, Russell CA,
228        Smith DJ, Rambaut A. 2014. Integrating influenza antigenic dynamics with molecular
229        evolution. eLife [Internet] 3:e01914–e01926. Available from:
230        http://elifesciences.org/lookup/doi/10.7554/eLife.01914

231    Ben Bolker, Team RDC. 2020. bbmle: Tools for General Maximum Likelihood Estimation.

232     Available from: https://CRAN.R-project.org/package=bbmle

233     Carrat F, Vergu E, Ferguson NM, Lemaitre M, Cauchemez S, Leach S, Valleron A-J. 2008. Time
234          Lines of Infection and Disease in Human Influenza: A Review of Volunteer Challenge
235          Studies. Am J Epidemiol 167:775–785.

236     De Maio N, Worby CJ, Wilson DJ, Stoesser N. 2018. Bayesian reconstruction of transmission
237          within outbreaks using genomic variants.Koelle K, editor. PLoS Comput Biol 14:e1006117–
238          e1006123.

239     Debbink K, McCrone JT, Petrie JG, Truscon R, Johnson E, Mantlo EK, Monto AS, Lauring AS.
240          2017. Vaccination has minimal impact on the intrahost diversity of H3N2 influenza
241          viruses.Perez DR, editor. PLoS Pathog 13:e1006194.

242     Didelot X, Fraser C, Gardy J, Colijn C. 2017. Genomic infectious disease epidemiology in
243          partially sampled and ongoing outbreaks. Molecular Biology and Evolution:msw075–11.

244     Didelot X, Gardy J, Colijn C. 2014. Bayesian Inference of Infectious Disease Transmission from
245          Whole-Genome Sequence Data. Molecular Biology and Evolution 31:1869–1879.

246     Geoghegan JL, Senior AM, Holmes EC. 2016. Pathogen population bottlenecks and adaptive
247          landscapes: overcoming the barriers to disease emergence. Proc. Biol. Sci.
248          283:20160727–20160729.

249     Hall M, Woolhouse M, Rambaut A. 2015. Epidemic Reconstruction in a Phylogenetics
250          Framework: Transmission Trees as Partitions of the Node Set.Salathé M, editor. PLoS
251          Comput Biol 11:e1004613–e1004636.

252     Kimura M. 1955. SOLUTION OF A PROCESS OF RANDOM GENETIC DRIFT WITH A
253          CONTINUOUS MODEL. Proceedings of the National Academy of Sciences 41:144–150.

254     Lakdawala SS, Jayaraman A, Halpin RA, Lamirande EW, Shih AR, Stockwell TB, Lin X,
255          Simenauer A, Hanson CT, Vogel L, et al. 2015. The soft palate is an important site of
256          adaptation for transmissible influenza viruses. Nature 526:122–125.

257     Lemey P, Suchard M, Rambaut A. 2009. Reconstructing the initial global spread of a human
258          influenza pandemic: A Bayesian spatial-temporal model for the global spread of H1N1pdm.
259          PLoS Curr 1:RRN1031.

260     Lumby CK, Zhao L, Breuer J, Illingworth CJR. 2020. A large effective population size for
261          established within-host influenza virus infection. eLife 9:217.

262     McCrone JT, Lauring AS. 2016. Measurements of intrahost viral diversity are extremely
263          sensitive to systematic errors in variant calling. J. Virol. 90:JVI.00667–16–6895.

264     McCrone JT, Woods RJ, Martin ET, Malosh RE, Monto AS, Lauring AS. 2018. Stochastic
265          processes constrain the within and between host evolution of influenza virus. eLife 7:24.

266     Rambaut A, Drummond AJ, Xie D, Baele G, Suchard MA. 2018. Posterior Summarization in
267          Bayesian Phylogenetics Using Tracer 1.7.Susko E, editor. Systematic Biology 67:901–904.

268   Rouzine IM, Rodrigo A, Coffin JM. 2001. Transition between stochastic evolution and
269       deterministic evolution in the presence of selection: general theory and application to
270       virology. Microbiology and Molecular Biology Reviews 65:151–185.

271   Sanjuán R, Nebot MR, Chirico N, Mansky LM, Belshaw R. 2010. Viral mutation rates. J. Virol.
272       84:9733–9748.

273   Sobel Leonard A, McClain MT, Smith GJD, Wentworth DE, Halpin RA, Lin X, Ransier A,
274       Stockwell TB, Das SR, Gilbert AS, et al. 2017. The effective rate of influenza reassortment
275       is limited during human infection. PLoS Pathog 13:e1006203.

276   Suchard MA, Lemey P, Baele G, Ayres DL, Drummond AJ, Rambaut A. 2018. Bayesian
277       phylogenetic and phylodynamic data integration using BEAST 1.10. Virus Evol 4:170.

278   Team RC. 2020. R: A Language and Environment for Statistical Computing. Available from:
279       https://www.R-project.org/

280   Visher E, Whitefield SE, McCrone JT, Fitzsimmons W, Lauring AS. 2016. The Mutational
281       Robustness of Influenza A Virus.Ferguson NM, editor. PLoS Pathog 12:e1005856.

282   Williamson EG, Slatkin M. 1999. Using maximum likelihood to estimate population size from
283       temporal changes in allele frequencies. Genetics 152:755–761.

284   Xue KS, Stevens-Ayers T, Campbell AP, Englund JA, Pergam SA, Boeckh M, Bloom JD. 2017.
285       Parallel evolution of influenza across multiple spatiotemporal scales. eLife 6:46.

286   Zhao L, Abbasi AB, Illingworth CJR. 2019. Mutational load causes stochastic evolutionary
287       outcomes in acute RNA viral infection. Virus Evol 5:686–12.

288

289

11

290   **Figure 1.** (A) Joint estimate of within-host mutation rate and effective population size. Contour

291   plot shows the log likelihood surface for estimates of the effective population size and neutral

292   mutation rate. The point represents the peak ($\mu = 3.4 \times 10^{-6}$, $N_e = 36$). Log likelihoods for each

293   contour are indicated. (B) The distribution of $N_e$ estimated in 500 subsamples of the data in

294   which one iSNV was taken per individual. The bimodality of the distribution reflects a slight

295   sensitivity to the inclusion of a few specific iSNV. (C) The posterior and prior probability

296   densities for $N_e$ over all values explored in the in the combined MCMC chains (22-93). The 95%

297   HPD of the posterior (32-72) is shaded blue.

298

299   **Table 1. Within host effective population size of IAV**
300

| iSNV Used | Generation Time (h) | Effective Population Size (95% CI) |
|---|---|---|
| All | 6 | 34 (25-46) |
| All | 12 | 17 (13-23) |
| Nonsynonymous | 6 | 27 (16-44) |
| Synonymous | 6 | 40 (27-59) |
| All | 12 | 17 (13-23) |
| Nonsynonymous | 12 | 14 (8-22) |
| Synonymous | 12 | 20 (14-29) |

301

Figure 1