

The evolutionary dynamics of influenza A virus within and between human hosts

Running Title: Influenza dynamics within and between hosts

John T. McCrone¹, Robert J. Woods², Emily T. Martin³, Ryan E. Malosh³, Arnold S. Monto³, and Adam S. Luring^{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

* Corresponding author

Adam S. Luring

1150 W. Medical Center Dr.

MSRB1 Room 5510B

Ann Arbor, MI 48109-5680

(734) 764-7731

aluring@med.umich.edu

Abstract Word Count:

Main Text Word Count:

Key Words – influenza virus, bottleneck, transmission, diversity, evolution

1 **Abstract**

2 A complete understanding of influenza virus evolution requires studies at all levels, as viral
3 evolutionary dynamics may differ across spatial and temporal scales. The relative contribution of
4 deterministic processes, such as selection, and stochastic processes, such as genetic drift, is
5 influenced by the virus' effective population size. While the global evolution of influenza A virus
6 (IAV) is dominated by the positive selection of novel antigenic variants that circulate in the
7 tropics, much less is known about the virus' evolution within and between human hosts. With
8 few exceptions, most of the available data derive from studies of chronically infected,
9 immunocompromised hosts, experimental infections with attenuated viruses, or animal models.
10 Here we define the evolutionary dynamics of IAV in human hosts through next generation
11 sequencing of 249 upper respiratory specimens from 200 individuals collected over 6290
12 person-seasons of observation. Because these viruses were collected over 5 seasons from
13 individuals in a prospective community-based cohort, they are broadly representative of natural
14 human infections with seasonal viruses. Within host genetic diversity was low, and we found
15 little evidence for positive selection of minority variants. We used viral sequence data from 35
16 serially sampled individuals to estimate a within host effective population size of 30-50. This
17 estimate is consistent across several models and robust to the models' underlying assumptions.
18 We also identified 43 epidemiologically linked and genetically validated transmission pairs.
19 Maximum likelihood optimization of multiple transmission models estimates an effective
20 transmission bottleneck of 1-2 distinct genomes. Our data suggest that positive selection of
21 novel viral variants is inefficient at the level of the individual host and that genetic drift and other
22 stochastic processes dominate the within and between host evolution of influenza A viruses.

23

24 **Introduction**

25 The rapid evolution of influenza viruses has led to reduced vaccine efficacy, widespread drug
26 resistance, and the continuing emergence of novel strains. Broadly speaking, evolution is the

27 product of deterministic processes, such as selection, and stochastic processes, such as
28 genetic drift (1). The relative contribution of each is greatly affected by the effective population
29 size, or size of an idealized population whose dynamics are similar to that of the population in
30 question (2). If the effective population size of a virus is large, as in quasispecies models,
31 evolution is largely deterministic and the frequency of a mutation can be predicted based on its
32 starting frequency and selection coefficient. In small populations, selection is inefficient, and
33 changes in mutation frequency are strongly influenced by genetic drift.

34

35 Viral dynamics may differ across spatial and temporal scales, and a complete understanding of
36 influenza evolution requires studies at all levels (3, 4). The global evolution of influenza A virus
37 (IAV) is dominated by the positive selection of novel antigenic variants that circulate in the
38 tropics and subsequently seed annual epidemics in the Northern and Southern hemisphere (5).
39 Whole genome sequencing has also demonstrated the importance of intrasubtype reassortment
40 to the emergence of diverse strains that differ in their antigenicity. While continual positive
41 selection of antigenically drifted variants drives global patterns, whole genome sequencing of
42 viruses on more local scales suggests the importance of stochastic processes such as strain
43 migration and within-clade reassortment (6).

44

45 With the advent of next generation sequencing, it is now feasible to efficiently sequence patient-
46 derived isolates at sufficient depth of coverage to define the diversity and dynamics of virus
47 evolution within individual hosts (7). Studies of IAV populations in animal and human systems
48 suggest that most intrahost single nucleotide variants (iSNV) are rare and that intrahost
49 populations are subject to strong purifying selection (8-14). While positive selection of adaptive
50 variants is commonly observed in cell culture (15-17), it has only been documented within
51 human hosts in the extreme cases of drug resistance (8, 18, 19), long-term infection of
52 immunocompromised hosts (20) or experimental infections with attenuated viruses (21). Indeed,

53 we and others have been unable to identify evidence for positive selection in natural human
54 infections (13, 14), and its relevance to within host processes is unclear.

55
56 Despite limited evidence for positive selection, it is clear that novel mutations do arise within
57 hosts. Their potential for subsequent spread through host populations is determined by the size
58 of the transmission bottleneck (22, 23). If the transmission bottleneck is sufficiently wide, low
59 frequency variants can plausibly be transmitted and spread through host populations (24).
60 Because the transmission bottleneck is conceptually similar to the effective population size
61 between hosts, its size will also inform the relative importance of selection and genetic drift in
62 determining which variants are transmitted. While experimental infections of guinea pigs and
63 ferrets suggest a very narrow transmission bottleneck (25, 26), studies of equine influenza
64 support a bottleneck wide enough to allow transmission of rare iSNV (9, 27). The only available
65 genetic study of influenza virus transmission in humans estimated a remarkably large
66 transmission bottleneck, allowing for transmission of 100-200 genomes (11, 28).

67
68 Here, we use next generation sequencing of within host influenza virus populations to define the
69 evolutionary dynamics of influenza A viruses (IAV) within and between human hosts. We apply
70 a benchmarked analysis pipeline to identify iSNV and to characterize the genetic diversity of
71 H3N2 and H1N1 populations collected over five post-pandemic seasons from individuals
72 enrolled in a prospective household study of influenza. We use these data to estimate the *in*
73 *vivo* mutation rate and the within and between host effective population size. We find that
74 intrahost populations are characterized by purifying selection, a small effective population size,
75 and limited positive selection. Contrary to what has been previously reported for human
76 influenza transmission (11), but consistent with what has been observed in other viruses (23),
77 we identify a very tight transmission bottleneck that limits the transmission of rare variants.

78

79 **Results**

80 We used next generation sequencing to characterize influenza virus populations collected from
81 individuals enrolled in the Household Influenza Vaccine Effectiveness (HIVE) study (29-32), a
82 community-based cohort that enrolls 213-340 households of 3 or more individuals in
83 Southeastern Michigan each year (Table 1). These households are followed prospectively from
84 October to April, with symptom-triggered collection of nasal and throat swab specimens for
85 identification of respiratory viruses by RT-PCR (see Methods). In contrast to case-ascertained
86 studies, which identify households based on an index case who seeks medical care, the HIVE
87 study identifies individuals regardless of illness severity. In the first four seasons of the study
88 (2010-2011 through 2013-2014), respiratory specimens were collected 0-7 days after illness
89 onset. Beginning in the 2014-2015 season, each individual provided two samples, a self-
90 collected specimen at the time of symptom onset and a clinic-collected specimen obtained 0-7
91 days later. Each year, 59-69% of individuals had self-reported or confirmed receipt of that
92 season's vaccine prior to local circulation of influenza virus.

93
94 Over five seasons and nearly 6,290 person-seasons of observation, we identified 77 cases of
95 influenza A/H1N1pdm09 infection and 313 cases of influenza A/H3N2 infection (Table 1).
96 Approximately half of the cases (n=166) were identified in the 2014-2015 season, in which there
97 was an antigenic mismatch between the vaccine and circulating strains (33). All other seasons
98 were antigenically matched. Individuals within a household were considered an
99 epidemiologically linked transmission pair if they were both positive for the same subtype of
100 influenza virus within 7 days of each other. Several households had 3 or 4 symptomatic cases
101 within this one-week window, suggestive of possible transmission chains (Table 1).

102

103 *Within host populations have low genetic diversity*

104 We processed all specimens for viral load quantification and next generation sequencing. Viral
105 load measurements (genome copies per μl) were used for quality control in variant calling,
106 which we have shown is highly sensitive to input titer (34) (Figure 1A). Accordingly, we report
107 data on 249 high quality specimens from 200 individuals, which had a viral load of $>10^3$ copies
108 per microliter of transport media, adequate RT-PCR amplification of all eight genomic segments,
109 and an average read coverage of $>10^3$ across the genome (Table 1, Supplemental Figure 1).

110
111 We identified intrahost single nucleotide variants (iSNV) using our empirically validated analysis
112 pipeline (34). Our approach relies heavily on the variant caller DeepSNV, which uses a clonal
113 plasmid control to distinguish between true iSNV and errors introduced during sample
114 preparation and/or sequencing (35). Given the diversity of influenza viruses that circulate locally
115 each season, there were a number of instances in which our patient-derived samples had
116 mutations that were essentially fixed (>0.95 frequency) relative to the clonal control. DeepSNV
117 is unable to estimate an error rate for the control or reference base at these positions. We
118 therefore performed an additional benchmarking experiment to identify a threshold for majority
119 iSNV at which we could correctly infer whether or not the corresponding minor allele was also
120 present (see Methods). We found that we could correctly identify a minor allele at a frequency of
121 $\geq 2\%$ when the frequency of the major allele was $\leq 98\%$. We therefore report data on iSNV
122 present at frequencies between 2 and 98%. As expected, this threshold improved the specificity
123 of our iSNV identification and decreased our sensitivity to detect variants below 5% compared to
124 our initial validation experiment (34), which did not employ a frequency threshold (Supplemental
125 Table 1).

126
127 Consistent with our previous studies and those of others, we found that the within host diversity
128 of human influenza A virus (IAV) populations is low (11, 13, 14, 21, 34). Two hundred forty-three
129 out of the 249 samples had fewer than 10 minority iSNV (median 2, IQR 1-3). There were 6

130 samples with greater than 10 minority iSNV. In 3 of these cases, the frequency of iSNVs were
131 tightly distributed about a mean suggesting that the iSNV were linked and that the samples
132 represented mixed infections. Consistent with this hypothesis, putative genomic haplotypes
133 based on these minority iSNV clustered with distinct isolates on phylogenetic trees
134 (Supplemental Figures 2 and 3). While viral shedding was well correlated with days post
135 symptom onset (Figure 1A) the number of minority iSNV identified was not affected by the day
136 of infection, viral load, subtype, or vaccination status (Figure 1B and Supplemental Figure 4).
137
138 The vast majority of minority variants were rare (frequency 0.02-0.07), and iSNV were
139 distributed evenly across the genome (Figure 1C and 1D). The ratio of nonsynonymous to
140 synonymous variants was 0.64 and was never greater than 1 in any 5% bin, which suggests
141 that within host populations were under purifying selection. We also found that minority variants
142 were rarely shared among multiple individuals. Ninety-five percent of minority iSNV were only
143 found once, 4.7% were found in 2 individuals, and no minority iSNV were found in more than 3
144 individuals. The low level of shared diversity suggests that within host populations were
145 exploring distinct regions of sequence space with little evidence for parallel evolution. Of the 31
146 minority iSNV that were found in multiple individuals (triangles in Figure 1D), 4 were
147 nonsynonymous.

148
149 Although the full range of the H3 antigenic sites have not been functionally defined, it is
150 estimated that 131 of the 329 amino acids in HA1 lie in or near these sites (36). We identified 17
151 minority nonsynonymous iSNV in these regions (Supplemental Table 2). Six of these were in
152 positions that differ among antigenically drifted viruses (37, 38), and two (193S and 189N) lie in
153 the “antigenic ridge” that is a major contributor to drift (39). Three of these have been detected
154 at the global level as consensus variants since the time of isolation (128A, 193S and 262N) with
155 two (193S and 262N) seemingly increasing in global frequency (40) (Supplemental Figure 5).

156 Additionally, we identified 1 putative H1N1 antigenic variant (208K in C_a) (41, 42). In total,
157 putative antigenic variants account for 1.0-2.5% of minority iSNV identified and were found in
158 3.5-7.5% of infections. None of these iSNV were shared among multiple individuals.

159

160 *Estimation of effective population size*

161 Given the above observations, we hypothesized that within host populations of IAV are under
162 purifying selection and that variants that rise to detectable levels do so by a neutral process as
163 opposed to positive selection. Consistent with this hypothesis, we found that nonsynonymous
164 and synonymous iSNV exhibited similar changes in frequency over time in the 35 individuals
165 who provided serial specimens that contained iSNV (Figure 2A and 2B). We used the diffusion
166 approximation to the Wright-Fisher model in conjunction with maximum likelihood estimation to
167 determine the within host effective population size (N_e) of IAV (43). This model assumes that
168 changes in iSNV frequency are due solely to random genetic drift and not selection, that iSNV
169 are independent of one another, and that the effective population is sufficiently large to justify a
170 continuous approximation to changes in allele frequency. While it is impossible to predict with
171 certainty the trajectory of allele frequencies under random genetic drift, the Wright-Fisher model,
172 and the diffusion approximation in particular, assigns probabilities to frequency changes given
173 an N_e and the number of generations between sample times. In our model we fixed the within
174 host generation time as either 6 or 12 hours (24) and report the findings for the 6 hour
175 generation time below. We then asked what population size makes the observed changes in
176 frequency most likely (Figure 2B). We restricted this analysis to samples taken at least 1 day
177 apart ($n = 29$), as there was very little change in iSNV frequency in populations sampled twice
178 on the same day ($R = 0.990$, Figure 2B and Supplemental Figure 6). The concordance of same
179 day samples suggests that our sampling procedure is reproducible and that less than a
180 generation had passed between samplings. Maximum likelihood optimization of this diffusion
181 model revealed a within host effective population size of 35 (95% CI 26-46, Table 2).

182

183 The diffusion approximation makes several simplifying assumptions, which if violated could
184 influence our findings. In particular, the model assumes a large population. To ensure our
185 results were robust to this assumption, we employed a discrete interpretation of the Wright-
186 Fisher model which makes no assumptions about population size (44). In this case we found an
187 effective population size of 32 (95% CI 28-41), very close to our original estimate (Table 2).
188 Both models assume complete independence of iSNV. To ensure this assumption did not affect
189 our results, we fit the discrete model 1000 times, each time randomly subsetting our data such
190 that only one iSNV per individual was included. This simulates a situation in which all modeled
191 iSNV are independent and our assumption is met. Under these conditions we found a median
192 effective population size of 33 (IQR 32-40), demonstrating negligible bias in the initial analysis
193 due to correlation between iSNV.

194

195 As above, most iSNV in the longitudinal samples were rare (< 10%) and many became extinct
196 between samplings. To ensure that our models were capable of accurately estimating the
197 effective population size from such data, we simulated 1000 Wright-Fisher populations with
198 iSNV present at approximately the same starting frequencies as in our data set an N_e of 30, 50,
199 or 100. In these simulations, we found mean N_e of 34, 56 and 117 (Figure 2C), which suggests
200 that our estimate is not an artifact of the underlying data structure.

201

202 To this point, we have assumed that neutral processes are responsible for the observed
203 changes in iSNV frequency within hosts. Although this assumption seems justified at least in
204 part by the analysis above, we tested the robustness of our models by fitting the
205 nonsynonymous ($n = 27$) and synonymous iSNV ($n = 36$) separately. Here, we estimated an
206 effective population size of 30 using the nonsynonymous iSNV and an effective population size
207 of 37 using the synonymous iSNV (Table 2). These estimates are very close to that derived

208 from the whole dataset and suggest that nonsynonymous and synonymous mutations are
209 influenced by similar within host processes. To further ensure that our results were not driven by
210 a few outliers subject to strong selection, we ranked iSNV by their change in frequency over
211 time and consecutively removed iSNV with the most extreme changes. We estimated the
212 effective population size at each iteration and found we would have to remove the top 75% most
213 extreme iSNV to increase the effective population size by a factor of 10 (Figure 2D). Therefore,
214 our estimates are robust to a reasonable number of non-neutral sites. Finally, we also applied a
215 separate Approximate Bayesian Computational (ABC) method, which uses a non-biased
216 moment estimator in conjunction with ABC to estimate the effective population size of a
217 population as well as selection coefficients for the iSNV present (17, 45). This distinct approach
218 relaxes our assumption regarding neutrality. We applied this analysis to the 16 longitudinal pairs
219 that were sampled 1 day apart and estimated an effective population of 54. We were unable to
220 reject neutrality for just 4 of the 35 iSNV in this data set (Figure 2E). These four mutations were
221 distributed between 2 individuals. Each individual had one nonsynonymous iSNV and one
222 synonymous iSNV. Neither were putative antigenic variants.

223

224 *Identification of forty-three transmission pairs*

225 The amount of diversity that passes between individuals during transmission determines the
226 extent to which within host evolution can affect larger evolutionary trends. We analyzed virus
227 populations from 85 households with concurrent infections to quantify the level of shared viral
228 diversity and to estimate the size of the IAV transmission bottleneck (Table 1). Because
229 epidemiological linkage does not guarantee that concurrent cases constitute a transmission pair
230 (46), we used a stringent rubric to eliminate individuals in a household with co-incident
231 community acquisition of distinct viruses. We considered all individuals in a household with
232 symptom onset within a 7-day window to be epidemiologically linked. The donor in each putative
233 pair was defined as the individual with the earlier onset of symptoms. We discarded a

234 transmission event if there were multiple possible donors with the same day of symptom onset.
235 Donor and recipients were not allowed to have symptom onset on the same day, unless the
236 individuals were both index cases for the household. In these 6 instances, we analyzed the data
237 for both possible donor-recipient directionalities. Based on these criteria, our cohort had 124
238 putative household transmission events over 5 seasons (Table 1). Of these, 52 pairs had
239 samples of sufficient quality for reliable identification of iSNV from both individuals.

240
241 We next used sequence data to determine which of these 52 epidemiologically linked pairs
242 represented true household transmission events as opposed to coincident community-acquired
243 infections. We measured the genetic distance between influenza populations from each
244 household pair by L1-norm and compared these distances to those of randomly assigned
245 community pairs within each season (Figure 3A, see also trees in Supplemental Figures 2 and
246 3). While the L1-norm of a pair captures differences between the populations at all levels, in our
247 cohort, it was largely driven by differences at the consensus level. We only considered
248 individuals to be a true transmission pair if they had a genetic distance below the 5th percentile
249 of the community distribution of randomly assigned pairs (Figure 3A). Forty-seven household
250 transmission events met this criterion (Figure 3B). Among these 47 sequence-validated
251 transmission pairs, 3 had no iSNV in the donor and 1 additional donor appeared to have a
252 mixed infection. These four transmission events were removed from our bottleneck analysis as
253 donors without iSNV are uninformative and mixed infections violate model assumptions of site
254 independence (see Methods). We estimated the transmission bottleneck in the remaining 43
255 high-quality pairs (37 H3N2, 6 H1N1, Figure 3B).

256
257 A transmission bottleneck restricts the amount of genetic diversity that is shared by both
258 members of a pair. We found that few minority iSNV were polymorphic in both the donor and
259 recipient populations (Figure 3C). Minority iSNV in the donor were either absent or fixed in the

260 recipient (top and bottom of plot). The lack of shared polymorphic sites (which would lie in the
261 middle of the plot in Figure 3C) suggests a stringent effective bottleneck in which only one allele
262 is passed from donor to recipient.

263

264 *Estimation of the transmission bottleneck*

265 We applied a simple presence-absence model to quantify the effective transmission bottleneck
266 in our cohort. True to its name, the presence-absence model simply measures whether or not a
267 donor allele is present or absent in the recipient sample. Under this model, transmission is a
268 neutral, random sampling process, and the probability of transmission is simply the probability
269 that the iSNV will be included at least once in the sample given its frequency in the donor and
270 the sample size, or bottleneck. We estimated a distinct bottleneck for each transmission pair
271 and assumed these bottlenecks followed a zero-truncated Poisson distribution. This model also
272 assumes that the sensitivity for detection of transmitted iSNVs is perfect and that each genomic
273 site is independent of all others. We then used maximum likelihood optimization to determine
274 the distribution of bottleneck sizes that best fit the data. We found a zero-truncated Poisson
275 distribution with a mean of 1.66 ($\lambda = 1.12$; 0.51-1.99, 95% CI) best described the data.
276 This distribution indicates that the majority of bottlenecks are 1, and that very few are greater
277 than 5 (probability 0.2%). There were no apparent differences between H3N2 and H1N1 pairs.
278 The model fit was evaluated by simulating each transmission event 1,000 times. The presence
279 or absence of each iSNV in the recipient was noted and the probability of transmission given
280 donor frequency determined. The range of simulated outcomes matched the data well, which
281 suggests that transmission is a selectively neutral event characterized by a stringent bottleneck
282 (Figure 3D).

283

284 The majority of transmitted iSNV were fixed in the recipients. Although this trend matches the
285 expectation given a small bottleneck, these data could also be consistent with a model in which

286 the probability of transmission is determined by the frequency at which iSNV are found at the
287 community level. To ensure our bottleneck estimates were an outcome of neutral transmission
288 and not an artifact of the larger community population structure or selection for the community
289 consensus, we created a null model by randomly assigning community “recipients” to each
290 donor in our transmission pairings. Each community “recipient” was drawn from the pool of
291 individuals that were infected after the donor but in the same season and with the same subtype
292 as the donor. We then identified whether or not each donor iSNV was found in the community
293 recipient and determined the relationship between donor frequency and probability of
294 “transmission” for 1,000 such simulations. Given the low level of diversity in our cohort, we
295 predicted that rare iSNV would be unlikely to be found in a random sample, while the major
296 alleles should be fixed in most random samples. This trend is clearly demonstrated in Figure 3E.
297 It is also clear that this null model fit the data much more poorly than the presence/absence
298 model, suggesting that the observed data in our bona fide transmission pairs were not a product
299 of community metapopulation structure, but rather an outcome of neutral sampling events.

300
301 Because our bottleneck estimates were much lower than what has previously been reported for
302 human influenza (11), we investigated the impact that our simplifying assumptions could have
303 on our results. In particular, the presence-absence model assumes perfect detection of variants
304 in donor and recipient, and it can therefore underestimate the size of a bottleneck in the setting
305 of donor-derived variants that are transmitted but not detected in the recipient. These “false
306 negative” variants can occur when the frequency of an iSNV drifts below the level of detection
307 (e.g. 2% frequency) or when the sensitivity of sequencing is less than perfect for variants at that
308 threshold (e.g. 15% sensitivity for variants at a frequency 2-5%). Leonard *et al.* recently
309 suggested that a beta binomial transmission model can account for the stochastic loss of
310 transmitted variants, by allowing for a limited amount of time-independent genetic drift within the
311 recipient (28). We modified this model to also account for our benchmarked sensitivity for rare

312 variants (Supplemental Table 1). For all donor-derived iSNV that were absent in the recipient,
313 we estimated the likelihood that these variants were transmitted but either drifted below our
314 level of detection or drifted below 10% and were missed by our variant identification. Despite the
315 relaxed assumptions provided by this modified beta binomial model, maximum likelihood
316 estimation only marginally increased the average bottleneck size (mean 1.71: lambda 1.19;
317 0.55-2.12, 95%CI) relative to the simpler presence-absence model. We simulated transmission
318 and subsequent random drift using the beta binomial model and the estimated bottleneck
319 distribution as above (Figure 3F). Although the model matched the data well, the fit was not
320 substantially better than that of the presence-absence model (AIC 75.5 for beta-binomial
321 compared to 76.7 for the presence-absence model).

322

323 *The mutation rate of influenza A virus within human hosts*

324 The stringent influenza transmission bottleneck suggests that most infections are founded by
325 one lineage and develop under essentially clonal processes. The diffusion approximation to the
326 Wright-Fisher model (see above and Figure 2) can be used to predict the rate at which
327 homogenous populations diversify from a clonal ancestor as a function of mutation rate and
328 effective population size (2). By applying maximum likelihood optimization to the model and the
329 frequency distribution of observed alleles (Figure 1C) we estimated an *in vivo* neutral mutation
330 rate of 4×10^{-6} mutations per nucleotide per replication cycle and a within host effective
331 population size of 33 (given a generation time of 6 hours). This is consistent with the estimates
332 above. As we have recently estimated that 13% of mutations in influenza A virus are neutral
333 (47), we estimated that the true *in vivo* mutation rate would be approximately 8 fold higher than
334 our neutral rate – on the order of $3-4 \times 10^{-5}$. This *in vivo* mutation rate is close to our recently
335 published estimate of influenza A mutation rates in epithelial cells by fluctuation test (48).

336

337 **Discussion**

338 We find that seasonal influenza A viruses replicate within and spread among human hosts with
339 very small effective population sizes. Because we used viruses collected over five influenza
340 seasons from individuals enrolled in a prospective household cohort, these dynamics are likely
341 to be broadly representative of many seasonal influenza infections. Other notable strengths of
342 our study include a validated sequence analysis pipeline and the use of models that are robust
343 to the underlying assumptions. The small effective size of intrahost populations and the tight
344 transmission bottleneck suggest that stochastic processes, such as genetic drift, dominate
345 influenza virus evolution at the level of individual hosts. This stands in contrast to prominent role
346 of positive selection in the global evolution of seasonal influenza.

347
348 While influenza virus populations are subject to continuous natural selection, selection is an
349 inefficient driver of evolution in small populations (2). Despite a large census, our findings
350 indicate that intrahost populations of influenza virus behave like much smaller populations. We
351 therefore expect random drift to be the major force driving the evolution of influenza virus within
352 human hosts. This finding contradicts previous studies, which have found signatures of adaptive
353 evolution in infected hosts (8, 19, 21, 49). However, these studies rely on data from infections in
354 which selective pressures are likely to be particularly strong (e.g. due to drug treatment or
355 infection with a poorly adapted virus), or in which the virus has been allowed to propagate for
356 extended periods of time. Under these conditions, one can identify the action of positive
357 selection on within host populations. We suggest that these are important exceptions to the drift
358 regime defined here.

359
360 We used both a simple presence-absence model and a more complex beta binomial model to
361 estimate an extremely tight transmission bottleneck. The small bottleneck size is driven by the
362 fact that within host diversity was low, and there were very few minority iSNV shared among
363 individuals in a transmission chain. While our methods for variant calling may be more

364 conservative than those used in similar studies, it is unlikely that our small bottleneck is an
365 artifact of this stringency. The beta binomial model accounts for false negative iSNV (i.e.
366 variants that are transmitted but not detected in the donor), which can lead to underestimated
367 transmission bottlenecks (28). Our formulation of this model incorporates empirically determined
368 sensitivity and specificity metrics to account for both false negative iSNV and false positive iSNV
369 (34). Furthermore, if rare, undetected, iSNV were shared between linked individuals, we would
370 expect to see transmission of more common iSNV (frequency 5-10%), which we can detect with
371 high sensitivity. In our dataset, however, the majority of minority iSNV above 5% were not
372 shared.

373
374 Although the size of our transmission bottleneck is consistent with estimates obtained for other
375 viruses and in experimental animal models of influenza (23, 25), it differs substantially from the
376 only other study of natural human infection (11, 28). While there are significant differences in the
377 design and demographics of the cohorts, the influenza seasons under study, and sequencing
378 methodology, the bottleneck size estimates are fundamentally driven by the amount of viral
379 diversity shared among individuals in a household. Importantly, we used both epidemiologic
380 linkage and the genetic relatedness of viruses in households to define transmission pairs and to
381 exclude confounding from the observed background diversity in the community. Whereas we
382 find that household transmission pairs and randomly assigned community pairs had distinct
383 patterns of shared consensus and minority variant diversity, Poon et al. found that rare iSNV
384 were often shared in both household pairs and randomly assigned community pairs (11).

385
386 Accurately modeling and predicting influenza virus evolution requires a thorough understanding
387 of the virus' population structure. Some models have assumed a large intrahost population and
388 a relatively loose transmission bottleneck (24, 50, 51). Here, adaptive iSNV can rapidly rise in
389 frequency and low frequency variants can have a high probability of transmission. In such a

390 model, it would be possible for the highly pathogenic H5N1 virus to develop the requisite 4-5
391 mutations to become transmissible through aerosols during a single acute infection of a human
392 host (50, 52). Although the dynamics of emergent avian influenza and human adapted seasonal
393 viruses likely differ, our work suggests that fixation of multiple mutations over the course of a
394 single acute infection is unlikely.

395
396 While it seems counterintuitive that influenza evolution is dominated by drift on local scales and
397 positive selection on global scales, these models are not necessarily in conflict. Within
398 individuals we have shown that the effective population is quite small, which suggests that
399 selection is inefficient. Indeed, we have deeply sequenced 332 intrahost populations from 283
400 individuals collected over more than 11,000 person-seasons of observation and only identified a
401 handful of minority antigenic variants with little evidence for positive selection (this work and
402 (14)). However, with several million infected individuals each year, even inefficient processes
403 and rare events are likely to happen at a reasonable frequency on a global scale.

404

405 **Methods**

406

407 *Description of the cohort*

408 The HIVE cohort (30, 31), established at the UM School of Public Health in 2010, enrolled and
409 followed households of at least 3 individuals with at least two children <18 years of age;
410 households were then followed prospectively throughout the year for ascertainment of acute
411 respiratory illnesses. Study participants were queried weekly about the onset of illnesses
412 meeting our standard case definition (two or more of: cough, fever/feverishness, nasal
413 congestion, sore throat, body aches, chills, headache if ≥ 3 yrs old; cough, fever/feverishness,
414 nasal congestion/runny nose, trouble breathing, fussiness/irritability, decreased appetite, fatigue
415 in <3 yrs old), and the symptomatic participants then attended a study visit at the research clinic

416 on site at UM School of Public Health for sample collection. For the 2010-2011 through 2013-
417 2014 seasons, a combined nasal and throat swab (or nasal swab only in children < 3 years of
418 age) was collected at the onsite research clinic by the study team. Beginning with the 2014-
419 2015 seasons, respiratory samples were collected at two time points in each participant meeting
420 the case definition; the first collection was a self- or parent-collected nasal swab collected at
421 illness onset. Subsequently, a combined nasal and throat swab (or nasal swab only in children <
422 3 years of age) was collected at the onsite research clinic by the study team. Families with very
423 young children (< 3 years of age) were followed using home visits by a trained medical assistant.

424

425 Active illness surveillance and sample collection for cases were conducted October through
426 May and fully captured the influenza season in Southeast Michigan in each of the study years.
427 Data on participant, family and household characteristics, and on high-risk conditions were
428 additionally collected by annual interview and review of each participant's electronic medical
429 record. In the current cohort, serum specimens were also collected twice yearly during fall
430 (November-December) and spring (May-June) for serologic testing for antibodies against
431 influenza.

432

433 This study was approved by the Institutional Review Board of the University of Michigan Medical
434 School, and all human subjects provided informed consent.

435

436 *Identification of influenza virus*

437 Respiratory specimens were processed daily to determine laboratory-confirmed influenza
438 infection. Viral RNA was extracted (Qiagen QIAamp Viral RNA Mini Kit) and tested by RT-PCR
439 for universal detection of influenza A and B. Samples with positive results by the universal
440 assay were then subtyped to determine A(H3N2), A(H1N1), A(pH1N1) subtypes and
441 B(Yamagata) and B(Victoria) lineages. We used primers, probes and amplification parameters

442 developed by the Centers for Disease Control and Prevention Influenza Division for use on the
443 ABI 7500 Fast Real-Time PCR System platform. An RNaseP detection step was run for each
444 specimen to confirm specimen quality and successful RNA extraction.

445

446 *Quantification of viral load*

447 Quantitative reverse transcription polymerase chain reaction (RT-qPCR) was performed on 5µl
448 RNA from each sample using CDC RT-PCR primers InfA Forward, InfA Reverse, and InfA
449 probe, which bind to a portion of the influenza M gene (CDC protocol, 28 April 2009). Each
450 reaction contained 5.4µl nuclease-free water, 0.5µl each primer/probe, 0.5µl SuperScript III
451 RT/Platinum Taq mix (Invitrogen 111732) 12.5µl PCR Master Mix, 0.1µl ROX, 5µl RNA. The
452 PCR master mix was thawed and stored at 4°C, 24 hours before reaction set-up. A standard
453 curve relating copy number to Ct value was generated based on 10-fold dilutions of a control
454 plasmid run in duplicate.

455

456 *Illumina library preparation and sequencing*

457 We amplified cDNA corresponding to all 8 genomic segments from 5µl of viral RNA using the
458 SuperScript III One-Step RT-PCR Platinum Taq HiFi Kit (Invitrogen 12574). Reactions consisted
459 of 0.5µl Superscript III Platinum Taq Mix, 12.5µl 2x reaction buffer, 6µl DEPC water, and 0.2µl
460 of 10µM Uni12/Inf1, 0.3µl of 10µM Uni12/Inf3, and 0.5µl of 10µM Uni13/Inf1 universal influenza
461 A primers (53). The thermocycler protocol was: 42°C for 60 min then 94°C for 2 min then 5
462 cycles of 94°C for 30 sec, 44°C for 30 sec, 68°C for 3 min, then 28 cycles of 94°C for 30 sec,
463 57°C for 30 sec, 68°C for 3 min. Amplification of all 8 segments was confirmed by gel
464 electrophoresis, and 750ng of each cDNA mixture were sheared to an average size of 300 to
465 400bp using a Covaris S220 focused ultrasonicator. Sequencing libraries were prepared using
466 the NEBNext Ultra DNA library prep kit (NEB E7370L), Agencourt AMPure XP beads (Beckman
467 Coulter A63881), and NEBNext multiplex oligonucleotides for Illumina (NEB E7600S). The final

468 concentration of each barcoded library was determined by Quanti PicoGreen dsDNA
469 quantification (ThermoFisher Scientific), and equal nanomolar concentrations were pooled.
470 Residual primer dimers were removed by gel isolation of a 300-500bp band, which was purified
471 using a GeneJet Gel Extraction Kit (ThermoFisher Scientific). Purified library pools were
472 sequenced on an Illumina HiSeq 2500 with 2x125 nucleotide paired end reads. All raw
473 sequence data have been deposited at the NCBI sequence read archive (BioProject submission
474 ID: SUB2951236). PCR amplicons derived from an equimolar mixture of eight clonal plasmids,
475 each containing a genomic segment of the circulating strain were processed in similar fashion
476 and sequenced on the same HiSeq flow cell as the appropriate patient derived samples. These
477 clonally derived samples served as internal controls to improve the accuracy of variant
478 identification and control for batch effects that confound sequencing experiments.

479

480 *Identification of iSNV*

481 Intrahost single nucleotide variants were identified in samples that had greater than 10^3
482 genomes/ μ l and an average coverage $>1000x$ across the genome. Variants were identified
483 using DeepSNV and scripts available at https://github.com/lauringlab/variant_pipeline as
484 described previously (34) with a few minor and necessary modifications. Briefly, reads were
485 aligned to the reference sequence (H3N2 2010-2011 & 2011-2012 : GenBank CY121496-503,
486 H3N2 2012-2013:GenBank KJ942680-8, H3N2 2014-2015 : Genbank CY207731-8, H1N1
487 GenBank : CY121680-8) using Bowtie2 (54). Duplicate reads were then marked and removed
488 using Picard (<http://broadinstitute.github.io/picard/>). We identified putative iSNV using DeepSNV.
489 Bases with phred <30 were masked. Minority iSNV (frequency $<50\%$) were then filtered for
490 quality using our empirically determined quality thresholds (p-value <0.01 DeepSNV, average
491 mapping quality >30 , average Phred >35 , average read position between 31 and 94). To control
492 for PCR errors in samples with lower input titers, all isolates with titers between 10^3 and 10^5
493 genomes/ μ l were processed and sequenced in duplicate. Only iSNV that were found in both

494 replicates were included in down stream analysis. The frequency of the variant in the replicate
495 with higher coverage at the iSNV location was assigned as the frequency of the iSNV. Finally,
496 any SNV with a frequency below 2% was discarded.

497
498 Given the diversity of the circulating strain in a given season, there were a number of cases in
499 which isolates contained mutations that were essentially fixed (>95%) relative to the plasmid
500 control. Often in these cases, the minor allele in the sample matched the major allele in the
501 plasmid control. We were, therefore, unable to use DeepSNV in estimating the base specific
502 error rate at this site for these minor alleles and required an alternative means of eliminating
503 true and false minority iSNV. To this end we applied stringent quality thresholds to these
504 putative iSNV and implemented a 2% frequency threshold. In order to ensure we were not
505 introducing a large number of false positive iSNV into our analysis, we performed the following
506 experiment. Perth (H3N2) samples were sequenced on the same flow cell as both the Perth and
507 Victoria (H3N2) plasmid controls. Minority iSNV were identified using both plasmid controls. This
508 allowed us to identify rare iSNV at positions in which the plasmid controls differed both with and
509 without the error rates provided by DeepSNV. We found that at a frequency threshold of 2% the
510 methods were nearly identical (NPV of 1, and PPV of 0.94 compared to DeepSNV).

511

512 *Overview of models for effective population size*

513 We estimated the effective population size using two separate interpretations of a Wright-Fisher
514 population. At its base, the Wright-Fisher model describes the expected changes in allele
515 frequency of an ideal population, which is characterized by non-overlapping generations, no
516 migration, no novel mutation, and no population structure. We then asked what size effective
517 population would make the changes in frequency observed in our dataset most likely. We
518 calculated these values using two applications of the Wright-Fisher model (i) a diffusion

519 approximation (43) and (ii) a maximum likelihood approach based on the discrete interpretation
520 (44).

521

522 For these estimates we restricted our analysis to longitudinal samples from a single individual
523 that were separated by at least 1 day and only used sites that were polymorphic in the initial
524 sample (29 of the 49 total serial sample pairs). We modeled only the iSNV that were the minor
525 allele at the first time point, and we assumed a within host generation time of either 6 or 12
526 hours as proposed by Geoghegan *et. al* (24).

527

528 *Diffusion approximation*

529 The diffusion approximation was first solved by Kimura in 1955 (43). This approximation to the
530 discrete Wright-Fisher model has enjoyed widespread use in population genetics as it allows
531 one to treat the random time dependent probability distribution of final allele frequencies as a
532 continuous function (e.g. (55-60)). Here, we also included the limitations in our sensitivity to
533 detect rare iSNV by integrating over regions of this probability density that were either below our
534 limit of detection or within ranges where we expect less than perfect sensitivity. Our adaptation
535 of Kimura's original work is below.

536

537 Let $P(p_0, p_t, t | N_e)$ be the time dependent probability of a variant drifting from an initial
538 frequency of p_0 to p_t over the course of t generations given an effective population size of N_e
539 where $0 < p_t < 1$.

540

541 The time dependent derivative of this probability has been defined using the forward
542 Kolmogorov equation and the solution is here adapted from Kimura, 1955 (43).

543

$$P(p_0, p_t, t | N_e) = \sum_{i=1}^{\infty} p_0 q_0 i(i+1)(2i+1)F(1-i, i+2, 2, p) \times F(1-i, i+2, 2, p_t) e^{-\left[\frac{i(i+1)}{2N_e}\right]t}$$

544 (1)

545

546 Where $q = 1 - p$ and F is the hypergeometric function. We approximated the infinite sum by
 547 summing over the first 50 terms. When we added an additional 50 terms (100 in total) we found
 548 no appreciable change in the final log likelihoods.

549

550 We denote the event that an allele is not observed at the second time point as $p_t \approx 0$ and the
 551 probability of such an event as $P(p_0, p_t \approx 0, t | N_e)$. This probability is given in equation 2 as the
 552 sum of the probability that the variant is lost by generation t (i.e. the other allele is fixed
 553 $P(q_0, 1, t | N_e)$), the probability that it is not detected due to the limit of detection (i.e. $P(p, p_t \approx$
 554 $0, t | 0 < p_t < 0.02, N_e)$) and the probability the variant is not detected due to low sensitivity for
 555 rare variant detection (i.e. $P(p_0, p_t \approx 0, t | 0.02 < p_t < 0.1, N_e)$). The probability of not observing
 556 an allele at the second time is then

557

$$P(p_0, p_t \approx 0, t | N_e) = P(q_0, 1, t | N_e) + P(p, p_t \approx 0, t | 0 < p_t < 0.02, N_e) + P(p_0, p_t \approx 0, t | 0.02 < p_t < 0.1, N_e)$$

558

559 (2)

560

561 The first term in equation 2 is adapted from Kimura, 1955 as

562

$$P(q_0, 1, t | N_e) = q_0 + \sum_{i=1}^{\infty} (2i+1)p_0q_0(-1)^i F(1-i, i+2, 2, q_0) e^{-[i(i+1)/2N_e]t}$$

563 (3)

564 Where q is defined as above. (Note that this is simply the probability of fixation for a variant at
565 initial frequency q). As in equation 1 the infinite sum was approximated with a partial sum of 50
566 terms.

567
568 The probability of the allele drifting below our limit of detection can be found by integrating
569 equation 1 between 0 and our limit of detection, 0.02. This was done numerically using the
570 python package `scipy` (61).

571

$$P(p, p_t \approx 0, t \mid 0 < p_t < 0.02, N_e) = \int_0^{0.02} P(p_0, p_t, t \mid N_e) dp_t \quad (4)$$

572

573 Finally, the probability of an iSNV being present at the second time point, but escaping detection,
574 is given by the integral of equation 1 between our benchmarked frequencies (0.02,0.05) times
575 the false negative rate for that range. Here, we assumed the entire range had the same
576 sensitivity as the benchmarked frequency at the lower bound and rounded recipient titers down
577 to the nearest \log_{10} titer (e.g. $10^3, 10^4, 10^5$). We also assumed perfect sensitivity above 10%.

578

$$P(p_0, p_t \approx 0, t \mid 0.02 < p_t < 0.1, N_e) = \sum_{f_i}^{[0.02, 0.05, 0.10]} (\text{FNR} \mid \text{Titer}_r, f_i) \int_{f_i}^{f_{i+1}} P(p_0, p_t, t \mid N_e) dp_t \quad (5)$$

579

580 Where $(\text{FNR} \mid \text{Titer}_r, f_i)$ is the false negative rate given the frequency and the sample titer (See
581 Supplemental Table 1) and $P(p_0, p_t, t \mid N_e)$ is defined in equation 1.

582

583 The log likelihood of a given population size is then simply the sum of the log of $P(p_0, p_t, t \mid N_e)$
584 for each minor allele in the data set, where either the position is polymorphic at time t (i.e.
585 equation 1) or the allele is observed as lost at time t (i.e. equation 2)

586

587 *Discrete Wright-Fisher estimation of N_e*

588 The diffusion approximation treats changes in frequency as a continuous process because it
589 assumes sufficiently large N_e . That assumption can be relaxed, and the effective population size
590 can be determined, by applying a maximum likelihood method developed by Williamsom and
591 Slaktin 1999 (44). In this model, the true allele frequencies move between discrete states (i.e.
592 the frequency must be of the form i/N_e where i is a whole number in the range $[0, N_e]$. In the
593 original application, allele counts were used, and sampling error was added to the model as a
594 binomial distribution with n determined by the sample size. Here, we use the frequencies
595 available from next generation sequencing and estimate sampling error as a normal distribution
596 with mean equal to the observed frequency and a standard deviation equal to that observed in
597 our benchmarking study for the 10^4 genomes/ μ l samples ($\sigma = 0.014$) (34).

598

599 In this model, the probability of observing an allele frequency shift from p_0 to p_t in t generations
600 provided an effective population of N_e is the probability of observing p_0 given some initial state
601 q_0 and the probability of the population having that state, times the probability of observing p_t
602 given some final state q_t and the probability of moving from the initial to the final state summed
603 across all possible states.

$$P(p_0, p_t | N_e) = \sum_{q_0, q_t} P(p_0 | q_0)P(q_0 | N_e)P(p_t | q_t)P(q_t | q_0, N_e)$$

604 (6)

605

606 Where p are the observed probabilities and q are the real ones (of the form i/N_e discussed
607 above). The likelihood of observing a given frequency p_x given a defined state q_x is given by the
608 likelihood of drawing p_x from a normal distribution with mean q_x and standard deviation 0.014.

609

$$P(p_x | q_x) = \text{Norm}(q_x, 0.014)$$

610 (7)

611 As in Williamson and Slatkin 1999, we assume a uniform prior on the initial state. Because we
 612 know that our specificity is near perfect (Supplemental Table 1) and we restrict our analysis to
 613 only polymorphic sites, the probability of any initial state is given by

614

$$P(q_0 | N_e) = \frac{1}{N_e - 1}$$

615 (8)

616

617 and finally the probability of moving from one state to another in t generations is given by

618

$$P(q_t, q_0 | N_e) = v_0 M^t v_t$$

619 (9)

620 Where M is a square transmission matrix with $C = N_e + 1$ rows and columns. Where $m_{i,j}$ is the
 621 probability of going from the i th configuration to the j th or the probability of drawing $j - 1$ out of
 622 binomial distribution with mean $(i - 1)/N_e$ and a sample size N_e . v_0 is a row vector of initial
 623 frequencies q_0 with 100% chance of initial state q_0 , and v_t is column vector of the frequencies at
 624 time point t with 100% chance of the final state. In other words v_0 is a row vector of C states
 625 with 0 everywhere except in the i th position where $\frac{i-1}{N_e} = q_0$, and v_t is a column vector of C
 626 states with 0 everywhere except the j th position where $\frac{j-1}{N_e} = q_t$

627 Using the scalar and cumulative properties of matrix multiplication equation 6 reduces to

628

$$P(p_0, p_t | N_e) = [0, P(p_0 | q_{0_2})P(q_{0_2} | N_e), \dots, P(p_0 | q_{0_{N_e-1}})P(q_{0_{N_e-1}} | N_e), 0] M^t \begin{bmatrix} P(p_t | q_{t_1}) \\ \vdots \\ P(p_t | q_{t_{N_e}}) \end{bmatrix}$$

629 (10)

630 The first and last entries in v_0 are 0 because we assume all measured sites represent
631 polymorphisms at the first time of sampling. As above, the log likelihood of a given population
632 size is then simply the sum of the log of $P(p_0, p_t, t | N_e)$ for each minor allele in the data set.

633

634 *Simulations*

635 To simulate within host evolution we set N_e in equation 10 to either 30, 50 or 100. For each
636 minor allele we used the closest available non-zero state given the effective population size as
637 the starting state. We then calculated the probability of moving to any other state and selected a
638 final state from this distribution. We then drew a final measured frequency from the normal
639 distribution accounting for measurement errors.

640

641 *ABC model*

642 We estimated both the effective population size and selection coefficients using the approximate
643 Bayesian computation (ABC) described in (17, 45) with the scripts provided in (45). In its current
644 implementation, this analysis requires the same time points for each sample, and we restricted
645 this analysis to longitudinal samples taken 1 day apart. This subset constitutes 16 of the 29
646 modeled longitudinal samples. Briefly, we subsampled polymorphic sites to 1,000x coverage to
647 estimate allele counts from frequency data as in (17). We then estimated the prior distribution of
648 the effective population size using 10,000 bootstrap replicates. We selected a uniform
649 distribution on the range [-0.5,0.5] as the prior distribution for the selection coefficients. The
650 posterior distributions were determined from accepting the top 0.01% of 100,000 simulations.

651

652 *Overview of models used for estimating the transmission bottleneck*

653 We model transmission as a simple binomial sampling process (28). In our first model, we
654 assume any transmitted iSNV, no matter the frequency, will be detected in the recipient. In the
655 second, we relax this assumption and account for false negative iSNV in the recipient. To

656 include the variance in the transmission bottlenecks between pairs we use maximum likelihood
657 optimization to fit the average bottleneck size assuming the distribution follows a zero-truncated
658 Poisson distribution.

659

660 *Presence/Absence model*

661 The presence/absence model makes many simplifying assumptions. We assume perfect
662 detection of all transmitted iSNV in the recipient. For each donor iSNV, we measure only
663 whether or not the variant is present in the recipient. Any iSNV that is not found in the recipient
664 is assumed to have not been transmitted. We also assume the probability of transmission is
665 determined only by the frequency of the iSNV in the donor at the time of sampling (regardless of
666 how much time passes between sampling and transmission). The probability of transmission is
667 simply the probability that the iSNV is included at least once in a sample size equal to the
668 bottleneck. Finally, we assume all genomic sites are independent of one another. For this
669 reason, we discarded the one case where the donor was likely infected by two strains as the
670 iSNV were certainly linked.

671

672 In our within host models, we only tracked minor alleles as in our data set we only ever find 2
673 alleles at each polymorphic site. In this case, the frequency of the major allele is simply one
674 minus the frequency of the minor allele. Because the presence/absence model is unaware of
675 the frequency of alleles in the recipient we must track both alleles at each donor polymorphic
676 site.

677

678 Let A_1 and A_2 be alleles in some donor j at some genomic site i . Let $P(A_x)$ be the probability
679 that the x allele is the only transmitted allele. There are then three possible outcomes for each
680 site. Either only A_1 is transmitted, only A_2 is transmitted, or both A_1 and A_2 are transmitted. The
681 probability of only one allele being transmitted given a bottleneck size of N_b is

682

$$P_{i,j}(A_x | N_b) = p_x^{N_b}$$

683 (11)

684 where p_x is the frequency of the x allele in the donor. In other words, this is simply the
685 probability of only drawing A_x in N_b draws.

686

687 The probability of both alleles being transmitted is given by

688

$$P_{i,j}(A_1, A_2 | N_b) = 1 - (p_1^{N_b} + p_2^{N_b})$$

689 (12)

690 where p_1 and p_2 are the frequencies of the alleles respectively. This is simply the probability of
691 not picking only A_1 or only A_2 in N_b draws.

692

693 This system could easily be extended to cases where there are more than 2 alleles present at a
694 site; however, that never occurs in our data set.

695

696 For ease we will denote the likelihood of observing the data at a polymorphic site i in each
697 donor j given the bottleneck size N_b as $P_{i,j}(N_b)$ where $P_{i,j}(N_b) = P_{i,j}(A_x | N_b)$ if only one allele is
698 transmitted and $P_{i,j}(N_b) = P_{i,j}(A_1, A_2 | N_b)$ if two alleles are transmitted.

699

700 The log likelihood of a bottleneck of size N_b is given by

701

$$LL(N_b) = \sum_j \sum_i \text{Ln}(P_{i,j})$$

702 (13)

703

704 where Ln is the natural log, and i, j refers to the i th polymorphic site in the j th donor. This is the
705 log of the probability of observing the data summed over all polymorphic sites across all donors.
706 Because the bottleneck size is likely to vary between individuals, we used maximum likelihood
707 to fit the bottleneck distribution as oppose to fitting a single bottleneck value. Under this model
708 we assumed the bottlenecks were distributed according to a zero-truncated Poisson distribution
709 parameterized by λ . The likelihood of observing the data given a polymorphic site i in donor j
710 and λ is

$$P_{i,j}(\lambda) = \sum_{N_b=1}^{\infty} P_{i,j}(N_b)P(N_b | \lambda)$$

711 (14)

712
713 where $P_{i,j}(N_b)$ is defined as above, $P(N_b | \lambda)$ is the probability of drawing a bottleneck of size N_b
714 from a zero-truncated Poisson distribution with a mean of $\frac{\lambda}{1-e^{-\lambda}}$. The sum is across all possible
715 N_b defined on $[1, \infty]$. For practical purposes, we only investigated bottleneck sizes up to 100, as
716 λ is quite small and the probability of drawing a bottleneck size of 100 from a zero-truncated
717 Poisson distribution with $\lambda = 10$ is negligible. We follow this convention whenever this sum
718 appears.

719
720 The log likelihood of λ for the data set is given by

721

$$LL(\lambda) = \sum_j \sum_i \text{Ln} \left(\sum_{N_b=1}^{\infty} P_{i,j}(N_b)P(N_b | \lambda) \right)$$

722 (15)

723
724 *Beta Binomial model*

725 The Beta binomial model is explained in detail in Leonard *et al.* (28). It is similar to the
726 presence/absence model in that transmission is modeled as a simple sampling process;
727 however, it loosens a few restricting assumptions. In this model, the frequencies of transmitted
728 variants are allowed to change between transmission and sampling according a beta distribution.
729 The distribution is not dependent on the amount of time that passes between transmission and
730 sampling, but rather depends on the size of the founding population (here assumed to equal to
731 N_b) and the number of variant genomes present in founding population k . Note the frequency in
732 the donor is assumed to be the same between sampling and transmission.

733
734 The equations below are very similar to those presented by Leonard *et al.* with two exceptions.
735 First, we fit a distribution to the bottleneck sizes in our cohort instead of fitting a single value,
736 and second because we know the sensitivity of our method to detect rare variants based on the
737 expected frequency and the titer, we can include the possibility that iSNV are transmitted but
738 are missed due to poor sensitivity. Because the beta binomial model is aware of the frequency
739 of the iSNV in the recipient, no information is added by tracking both alleles at a genomic site i .
740 Let p_{i,j_d} represent the frequency of the minor allele frequency at position i in the donor of some
741 transmission pair j . Similarly, let p_{i,j_r} be the frequency of that same allele in the recipient of the
742 j th transmission pair. Then, as in Leonard *et al.*, the likelihood of some bottleneck N_b for the
743 data at site i in pair j where the minor allele is transmitted is given by

$$L(N_b)_{i,j} = \sum_{k=1}^{N_b} p_beta(p_{i,j_r} | k, N_b - k) p_bin(k | N_b, p_{i,j_d})$$

745 (16)

746
747 Where p_beta is the probability density function for the beta distribution and p_bin is the
748 probability mass function for the binomial distribution.

749

750 This is the probability density that the transmitted allele is found in the recipient at a frequency
751 of p_{i,j_r} given that the variant was in k genomes in a founding population of size N_b times the
752 probability k variant genomes would be drawn in a sample size of N_b from the donor where the
753 variant frequency was p_{i,j_d} . This is then summed for all possible k where $1 < k \leq N_b$.

754 As in equation 14 the likelihood of a zero truncated Poisson with a mean of $\frac{\lambda}{1-e^{-\lambda}}$ given this
755 transmitted variants is then given by

756

$$L(\lambda)_{i,j}^{\text{transmitted}} = \sum_{N_b=1}^{\infty} L(N_b)_{i,j} P(N_b | \lambda)$$

757 (17)

758

759 This is simply the likelihood of each N_b weighted by the probability of drawing a bottleneck size
760 of N_b from bottleneck distribution.

761 In this model, there are three possible mechanisms for a donor iSNV to not be detected in the
762 recipient. (i) The variant was not transmitted. (ii) The variant was transmitted but is present
763 below our level of detection (2%). (iii) The variant was transmitted and present above our level
764 of detection but represents a false negative in iSNV identification.

765

766 As in Leonard *et al.*, the likelihood of scenarios (i) and (ii) for a given N_b are expressed as

767

$$L(N_b)_{i,j}^{\text{lost}} = \sum_{k=0}^{N_b} \text{p_beta_cdf}(p_{i,j_r} < 0.02 | k, N_b - k) \text{p_bin}(k | N_b, p_{i,j_d})$$

768 (18)

769

770 Where p_beta_cdf is the cumulative distribution function for the beta distribution. Note that if
 771 $k = 0$ (i.e. the iSNV was not transmitted) then the term reduces to the probability of not drawing
 772 the variant in N_b draws.

773
 774 The likelihood of the variant being transmitted but not detected in the recipient given a
 775 bottleneck of N_b is described by

$$L(N_b)_{i,j}^{\text{missed}} = \sum_{k=0}^{N_b} \sum_{f_e \in [0.02, 0.05, 0.1]} p_beta_cdf(f_e < p_{i,j_r} < f_{e+1} | k, N_b - k) \times p_bin(k | N_b, p_{i,j_d})(FNR | Titer_r, f_e) \quad (19)$$

776
 777
 778
 779
 780 This is the likelihood of the variant existing in the ranges [0.02,0.05] or [0.05,0.1] given an initial
 781 frequency of k/N_b and a bottleneck size of N_b multiplied by the expected False Negative Rate
 782 (FNR) given the titer of the recipient and the lower frequency bound. As in our diffusion model,
 783 we assumed perfect sensitivity for detection of iSNV present above 10%, rounded recipient
 784 titers down to the nearest \log_{10} titer (e.g. $10^3, 10^4, 10^5$) and assumed the entire range $[f_e, f_{e+1}]$
 785 has the same sensitivity as the lower bound.

786
 787 The likelihood of λ for iSNV that are not observed in the recipient is then given by summing
 788 equations 18 and 19 across all possible N_b .

$$L(\lambda)_{i,j}^{\text{nontransmitted}} = \sum_{N_b=1}^{\infty} (L(N_b)_{i,j}^{\text{lost}} + L(N_b)_{i,j}^{\text{missed}}) P(N_b | \lambda) \quad (20)$$

790
 791

792 The log likelihood of the total dataset is then determined by summing log of equations 17 and 20
793 (as applicable) across all polymorphic sites in each donor. (As before here we sum of N_b within
794 the range [1,100].)

795

796 *Simulation*

797 In order evaluate the fits of the two transmission models, we simulated whether or not each
798 donor iSNV was transmitted or not. This involved converting each model to a presence absence
799 model. In each simulation, we assigned a bottleneck from the bottleneck distribution for each
800 transmission pair. We then determined the probability of only transmitting one allele (A_x where
801 $x \in [1,2]$ as in the presence/absence model above) and the probability of transmitted both
802 alleles (A_1, A_2 above) for each polymorphic site.

803

804 For the presence/absence model, the probabilities for each possible outcome are given by
805 equations 11 and 12. For the beta binomial model, the probability of only observing A_x at site i
806 is given by

$$P(A_x | N_b) = L(N_b)_{i,j}^{\text{lost}} + L(N_b)_{i,j}^{\text{missed}} \quad (21)$$

807

808

809 where $L(N_b)_{i,j}^{\text{lost}}$ and $L(N_b)_{i,j}^{\text{missed}}$ are defined as in equations 18 and 19 respectively, but with
810 $p_{i,j,d}$ replaced by $1 - p_{i,j,d}$. This is simply the probability of not observing the other allele in the
811 recipient.

812

813 Again, the probability of observing both alleles is

814

$$P(A_1, A_2 | N_b) = 1 - (P(A_1) + P(A_2))$$

815 (22)

816 where $P(A_1)$ and $P(A_2)$ are defined as in equation 21.

817

818 *Fitting mutation rate and N_e*

819 The diffusion approximation to the Wright - Fischer model allows us to make predictions on the
820 allele frequency spectrum of a population given a mutation rate and an effective population size.

821 The probability of observing a mutation at frequency p_t given an initial frequency of 0 can be
822 approximated as in (2)

823

$$P(0, p_t, t, | \mu, N_e) = \frac{2\mu N_e}{p_t} e^{-\frac{2N_e p_t}{t}}$$

824 (23)

825 Where μ is the mutation rate. In this model mutation increases an allele's frequency from 0 but
826 after that initial jump, drift is responsible for allowing the mutation to reach it's observed
827 frequency. Because the limit of equation 23 approaches infinity as p_t approaches 0 and for ease
828 in numerical integration, we assumed that any variant present at less than 0.1% was essentially
829 at 0%.

830

831 We then assumed each infection began as a clonal infection matching the consensus sequence
832 observed at the time of sampling. The likelihood of observing minor alleles at the observed
833 frequency is the given by equation 23.

834

835 As in the other within host models, we can account for nonpolymorphic sites by adding the
836 likelihood that no mutation is present $P(0, p_t \approx 0, t | p_t < 0.001, \mu, N_e)$, that a mutation is present
837 but below our level of detection $P(0, p_t \approx 0, t | p_t < 0.02, \mu, N_e)$, and that a mutation is present
838 but missed due to low sensitivity at low frequencies $P(0, p_t \approx 0, t | 0.02 < p_t < 0.1, \mu, N_e)$. In this

839 model we assumed 13133 mutagenic targets in each sample (the number of coding sites
840 present in the reference strain from 2014-2015).

841

842 The probability of not observing a mutation is given by

843

$$\begin{aligned}
 P(0, p_t \approx 0, t, | \mu, N_e) &= P(0, p_t \approx 0, t | p_t < 0.001, \mu, N_e) + \\
 &P(0, p_t \approx 0, t | p_t < 0.02, \mu, N_e) + \\
 &P(0, p_t \approx 0, t | 0.02 < p_t < 0.1, \mu, N_e)
 \end{aligned}
 \tag{24}$$

847 Where

$$P(0, p_t \approx 0, t | p_t < 0.001, \mu, N_e) = 1 - \int_{0.001}^1 P(0, p_t, t, | \mu, N_e) dp_t
 \tag{25}$$

848

849 and

$$P(0, p_t \approx 0, t | p_t < 0.02, \mu, N_e) = \int_{0.001}^{0.02} P(0, p_t, t, | \mu, N_e) dp_t
 \tag{26}$$

850

851 and

$$P(0, p_t \approx 0, t | 0.02 < p_t < 0.1, \mu, N_e) = \sum_{f_i}^{[0.02, 0.05, 0.10]} (\text{FNR} | \text{Titer}_r, f_i) \int_f^{f_{i+1}} P(p_0, p_t, t | \mu, N_e) dp_t
 \tag{27}$$

852

853

854 Where we follow the same convention as in equation 5 for determining the false negative rate.

855 The log likelihood of a given μ and N_e pair is then the sum of the log of equations 23 and 24 for

856 all possible sites in the data set.

857

858 Annotated computer code for all analyses can be accessed at

859 https://github.com/lauringlab/Host_level_IAV_evolution

860 **Acknowledgments**

861 This work was supported by a Clinician Scientist Development Award from the Doris Duke
862 Charitable Foundation (CSDA 2013105), a University of Michigan Discovery Grant, and R01
863 AI118886, all to ASL. The HIVE cohort was supported by NIH R01 AI097150 and CDC U01
864 IP00474 to ASM. JTM was supported by the Michigan Predoctoral Training Program in Genetics
865 (T32GM007544). RJW was supported by K08AI119182. We thank Alexey Kondrashov and
866 Aaron King for helpful discussion.

867

868 **References**

- 869 1. Kouyos RD, Althaus CL, Bonhoeffer S (2006) Stochastic or deterministic: what is the
870 effective population size of HIV-1? *Trends in Microbiology* 14(12):507–511.
- 871 2. Rouzine IM, Rodrigo A, Coffin JM (2001) Transition between stochastic evolution and
872 deterministic evolution in the presence of selection: general theory and application to
873 virology. *Microbiology and Molecular Biology Reviews* 65(1):151–185.
- 874 3. Nelson MI, Holmes EC (2007) The evolution of epidemic influenza. *Nat Rev Genet*
875 8(3):196–205.
- 876 4. Holmes EC (2009) RNA virus genomics: a world of possibilities. *Journal of Clinical*
877 *Investigation* 119(9):2488–2495.
- 878 5. Rambaut A, et al. (2008) The genomic and epidemiological dynamics of human influenza
879 A virus. *Nature* 453(7195):615–619.
- 880 6. Nelson MI, et al. (2006) Stochastic Processes Are Key Determinants of Short-Term
881 Evolution in Influenza A Virus. *PLoS Pathog* 2(12):e125.
- 882 7. Kao RR, Haydon DT, Lycett SJ, Murcia PR (2014) Supersize me: how whole-genome
883 sequencing and big data are transforming epidemiology. *Trends in Microbiology*:1–10.
- 884 8. Rogers MB, et al. (2015) Intrahost dynamics of antiviral resistance in influenza a virus
885 reflect complex patterns of segment linkage, reassortment, and natural selection. *mBio*
886 6(2):e02464–14.
- 887 9. Murcia PR, et al. (2010) Intra- and Interhost Evolutionary Dynamics of Equine Influenza
888 Virus. *J Virol* 84(14):6943–6954.
- 889 10. Iqbal M, et al. (2009) Within-host variation of avian influenza viruses. *Philosophical*
890 *Transactions of the Royal Society B: Biological Sciences* 364(1530):2739–2747.

- 891 11. Poon LLM, et al. (2016) Quantifying influenza virus diversity and transmission in humans.
892 *Nat Genet* 48(2):195–200.
- 893 12. Ghedin E, et al. (2009) Mixed Infection and the Genesis of Influenza Virus Diversity. *J*
894 *Viro* 83(17):8832–8841.
- 895 13. Dinis JM, et al. (2016) Deep Sequencing Reveals Potential Antigenic Variants at Low
896 Frequencies in Influenza A Virus-Infected Humans. *J Virol* 90(7):3355–3365.
- 897 14. Debbink K, et al. (2017) Vaccination has minimal impact on the intrahost diversity of
898 H3N2 influenza viruses. *PLoS Pathog* 13(1):e1006194.
- 899 15. Doud MB, Hensley SE, Bloom JD (2017) Complete mapping of viral escape from
900 neutralizing antibodies. *PLoS Pathog* 13(3):e1006271.
- 901 16. Archetti I, HORSFALL FL (1950) Persistent antigenic variation of influenza A viruses after
902 incomplete neutralization in ovo with heterologous immune serum. *The Journal of*
903 *Experimental Medicine* 92(5):441–462.
- 904 17. Foll M, et al. (2014) Influenza Virus Drug Resistance: A Time-Sampled Population
905 Genetics Perspective. *PLoS Genet* 10(2):e1004185–17.
- 906 18. Gubareva LV, Kaiser L, Matrosovich MN, Soo-Hoo Y, Hayden FG (2001) Selection of
907 Influenza Virus Mutants in Experimentally Infected Volunteers Treated with Oseltamivir.
908 *Journal of Infectious Diseases* 183(4):523–531.
- 909 19. Ghedin E, et al. (2010) Deep Sequencing Reveals Mixed Infection with 2009 Pandemic
910 Influenza A (H1N1) Virus Strains and the Emergence of Oseltamivir Resistance. *Journal*
911 *of Infectious Diseases* 203(2):168–174.
- 912 20. Xue KS, et al. (2017) Parallel evolution of influenza across multiple spatiotemporal
913 scales. *eLife* 6:46.
- 914 21. Sobel Leonard A, et al. (2016) Deep Sequencing of Influenza A Virus from a Human
915 Challenge Study Reveals a Selective Bottleneck and Only Limited Intrahost Genetic
916 Diversification. *J Virol* 90(24):11247–11258.
- 917 22. Alizon S, Luciani F, Regoes RR (2011) Epidemiological and clinical consequences of
918 within-host evolution. *Trends in Microbiology* 19(1):24–32.
- 919 23. Zwart MP, Elena SF (2015) Matters of Size: Genetic Bottlenecks in Virus Infection and
920 Their Potential Impact on Evolution. *Annual Review of Virology* 2(1):161–179.
- 921 24. Geoghegan JL, Senior AM, Holmes EC (2016) Pathogen population bottlenecks and
922 adaptive landscapes: overcoming the barriers to disease emergence. *Proc Biol Sci*
923 283(1837):20160727–9.
- 924 25. Varble A, et al. (2014) Influenza A Virus Transmission Bottlenecks Are Defined by
925 Infection Route and Recipient Host. *Cell Host Microbe* 16(5):691–700.
- 926 26. Wilker PR, et al. (2013) Selection on haemagglutinin imposes a bottleneck during

- 927 mammalian transmission of reassortant H5N1 influenza viruses. *Nat Comms* 4:1–11.
- 928 27. Hughes J, et al. (2012) Transmission of Equine Influenza Virus during an Outbreak Is
929 Characterized by Frequent Mixed Infections and Loose Transmission Bottlenecks. *PLoS*
930 *Pathog* 8(12):e1003081.
- 931 28. Sobel Leonard A, Weissman DB, Greenbaum B, Ghedin E, Koelle K (2017) Transmission
932 Bottleneck Size Estimation from Pathogen Deep-Sequencing Data, with an Application to
933 Human Influenza A Virus. *J Virol* 91(14):e00171–17.
- 934 29. Ohmit SE, et al. (2016) Substantial Influenza Vaccine Effectiveness in Households With
935 Children During the 2013–2014 Influenza Season, When 2009 Pandemic Influenza
936 A(H1N1) Virus Predominated. *Journal of Infectious Diseases* 213(8):1229–1236.
- 937 30. Ohmit SE, et al. (2015) Influenza Vaccine Effectiveness in Households With Children
938 During the 2012–2013 Season: Assessments of Prior Vaccination and Serologic
939 Susceptibility. *Journal of Infectious Diseases* 211(10):1519–1528.
- 940 31. Monto AS, Malosh RE, Petrie JG, Thompson MG, Ohmit SE (2014) Frequency of acute
941 respiratory illnesses and circulation of respiratory viruses in households with children over
942 3 surveillance seasons. *J Infect Dis* 210(11):1792–1799.
- 943 32. Petrie JG, et al. (2013) Influenza transmission in a cohort of households with children:
944 2010–2011. *PLoS ONE* 8(9):e75339.
- 945 33. Chung JR, et al. (2017) Prior season vaccination and risk of influenza during the 2014–
946 2015 season in the U.S. *J Infect Dis* 216(2):284–285.
- 947 34. McCrone JT, Lauring AS (2016) Measurements of intrahost viral diversity are extremely
948 sensitive to systematic errors in variant calling. *J Virol* 90(15):JVI.00667–16–6895.
- 949 35. Gerstung M, et al. (2012) Reliable detection of subclonal single-nucleotide variants in
950 tumour cell populations. *Nat Comms* 3:811.
- 951 36. Lee M-S, Chen JS-E (2004) Predicting Antigenic Variants of Influenza A/H3N2 Viruses.
952 *Emerg Infect Dis* 10(8):1385–1390.
- 953 37. Smith DJ, et al. (2004) Mapping the antigenic and genetic evolution of influenza virus.
954 *Science* 305(5682):371–376.
- 955 38. Wiley DC, Wilson IA, Skehel JJ (1981) Structural identification of the antibody-binding
956 sites of Hong Kong influenza haemagglutinin and their involvement in antigenic variation.
957 *Nature*.
- 958 39. Koel BF, Burke DF, Bestebroer TM (2013) Substitutions near the receptor binding site
959 determine major antigenic change during influenza virus evolution.
960 doi:10.1126/science.1240537.
- 961 40. Neher RA, Bedford T (2015) nextflu: real-time tracking of seasonal influenza virus
962 evolution in humans. *Bioinformatics* 31(21):3546–3548.

- 963 41. Caton AJ, Brownlee GG, Yewdell JW, Gerhard W (1982) The antigenic structure of the
964 influenza virus A/PR/8/34 hemagglutinin (H1 subtype). *Cell* 31(2 Pt 1):417–427.
- 965 42. Xu R, Ekiert DC, Krause JC, Hai R, Crowe JE (2010) Structural basis of preexisting
966 immunity to the 2009 H1N1 pandemic influenza virus. ... 328(5976):354–357.
- 967 43. Kimura M (1955) Solution of a process of random genetic drift with a continuous model.
- 968 44. Williamson EG, Slatkin M (1999) Using maximum likelihood to estimate population size
969 from temporal changes in allele frequencies. *Genetics* 152(2):755–761.
- 970 45. Foll M, Shim H, Jensen JD (2015) WFABC: a Wright-Fisher ABC-based approach for
971 inferring effective population sizes and selection coefficients from time-sampled data. *Mol*
972 *Ecol Resour* 15(1):87–98.
- 973 46. Petrie JG, et al. (2017) Application of an Individual-Based Transmission Hazard Model for
974 Estimation of Influenza Vaccine Effectiveness in a Household Cohort. *Am J Epidemiol*.
975 doi:10.1093/aje/kwx217.
- 976 47. Visher E, Whitefield SE, McCrone JT, Fitzsimmons W, Lauring AS (2016) The Mutational
977 Robustness of Influenza A Virus. *PLoS Pathog* 12(8):e1005856.
- 978 48. Pauly MD, Procaro MC, Lauring AS (2017) A novel twelve class fluctuation test reveals
979 higher than expected mutation rates for influenza A viruses. *eLife* 6:686.
- 980 49. Gubareva LV, Kaiser L, Matrosovich MN, Soo-Hoo Y, Hayden FG (2001) Selection of
981 Influenza Virus Mutants in Experimentally Infected Volunteers Treated with Oseltamivir.
982 *Journal of Infectious Diseases* 183(4):523–531.
- 983 50. Russell CA, et al. (2012) The Potential for Respiratory Droplet-Transmissible A/H5N1
984 Influenza Virus to Evolve in a Mammalian Host. *Science* 336(6088):1541–1547.
- 985 51. Peck KM, Chan CHS, Tanaka MM (2015) Connecting within-host dynamics to the rate of
986 viral molecular evolution. *Virus Evol* 1(1). doi:10.1093/ve/vev013.
- 987 52. Herfst S, et al. (2012) Airborne Transmission of Influenza A/H5N1 Virus Between Ferrets.
988 *Science* 336(6088):1534–1541.
- 989 53. Zhou B, et al. (2009) Single-reaction genomic amplification accelerates sequencing and
990 vaccine production for classical and Swine origin human influenza a viruses. *J Virol*
991 83(19):10309–10313.
- 992 54. Langmead B, Salzberg SL (2012) Fast gapped-read alignment with Bowtie 2. *Nat*
993 *Methods* 9(4):357–359.
- 994 55. Zanini F, Puller V, Brodin J, Albert J, Neher RA (2017) In vivo mutation rates and the
995 landscape of fitness costs of HIV-1. *Virus Evol* 3(1):1–12.
- 996 56. Myers S, Fefferman C, Patterson N (2008) Can one learn history from the allelic
997 spectrum? *Theoretical Population Biology* 73(3):342–348.
- 998 57. Nagasawa M, Maruyama T (1979) An application of time reversal of Markov processes to

- 999 a problem of population genetics. *Advances in Applied Probability*.
- 1000 58. Kimura M (1971) Theoretical foundation of population genetics at the molecular level.
1001 *Theoretical Population Biology* 2(2):174–208.
- 1002 59. Kimura M, Ohta T (1969) The Average Number of Generations until Fixation of a Mutant
1003 Gene in a Finite Population. *Genetics* 61(3):763–771.
- 1004 60. Kimura M (1964) Diffusion Models in Population Genetics. *Journal of Applied Probability*
1005 1(2):177.
- 1006 61. Oliphant TE (2007) Python for Scientific Computing. *Comput Sci Eng* 9(3):10–20.
- 1007
- 1008

1009 **Figure Legends**

1010 **Figure 1.** Within host diversity of IAV populations. (A) Boxplots (median, 25th and 75th
1011 percentiles, whiskers extend to most extreme point within median \pm 1.5 x IQR) of the number of
1012 viral genomes per microliter transport media stratified by day post symptom onset. Notches
1013 represent the approximate 95% confidence interval of the median. (B) Boxplots (median, 25th
1014 and 75th percentiles, whiskers extend to most extreme point within median \pm 1.5 x IQR) of the
1015 number of iSNV in 249 high quality samples stratified by day post symptom onset. (C)
1016 Histogram of within host iSNV frequency in 249 high quality samples. Bin width is 0.05
1017 beginning at 0.02. Mutations are colored nonsynonymous (blue) and synonymous (gold) (D)
1018 Location of all identified iSNV in the influenza A genome. Mutations are colored
1019 nonsynonymous (blue) and synonymous (gold) relative to that sample's consensus sequence.
1020 Triangles signify mutations that were found in more than one individual in a given season.

1021
1022 **Figure 2.** Within host dynamics of IAV. (A) Timing of sample collection for 35 paired longitudinal
1023 samples relative to day of symptom onset. Of the 49 total, 35 pairs had minor iSNV present in
1024 the first sample. (B) The change in frequency over time for minority nonsynonymous (blue) and
1025 synonymous (gold) iSNV identified for the paired samples in (A). (C) The distribution of effective
1026 population sizes estimated from 1,000 simulated populations. Simulations were run on
1027 populations with characteristics similar to the actual patient-derived populations and with the
1028 specified effective population size (x-axis). (D) The effect of iteratively removing iSNV with the
1029 most extreme change in frequency (fraction of iSNV removed, x-axis) on the estimated effective
1030 population size. The point represents the estimate when all iSNV are included (32). (E) The
1031 posterior distributions of selection coefficients estimated for the 35 iSNV present in isolates
1032 sampled one day apart. Distributions are colored according to class relative to the sample
1033 consensus sequence, nonsynonymous (blue) synonymous (gold). Variants for which the 95%

1034 highest posterior density intervals exclude 0.0 are noted in the margin.

1035

1036 **Figure 3.** Between host dynamics of IAV. (A) The distribution of pairwise L1-norm distances for
1037 household (blue) and randomly-assigned community (gold) pairs. The bar heights are
1038 normalized to the height of the highest bar for each given subset (47 for household, 1,590 for
1039 community). The red line represents the 5th percentile of the community distribution. (B) Timing
1040 of symptom onset for 52 epidemiologically linked transmission pairs. Day of symptom onset for
1041 both donor and recipient individuals is indicated by black dots. Dashed lines represent pairs that
1042 were removed due to abnormally high genetic distance between isolates, see (A). (C) The
1043 frequency of donor iSNV in both donor and recipient samples. Frequencies below 2% and
1044 above 98% were set to 0% and 100% respectively. (D) The presence-absence model fit
1045 compared with the observed data. The x-axis represents the frequency of donor iSNV with
1046 transmitted iSNV plotted along the top and nontransmitted iSNV plotted along the bottom. The
1047 black line indicates the probability of transmission for a given iSNV frequency as determined by
1048 logistic regression. Similar fits were calculated for 1,000 simulations with a mean bottleneck size
1049 of 1.66. Fifty percent of simulated outcomes lie in the darkly shaded region and 95% lie in the
1050 lightly shaded regions. (E) The outcome from 1,000 simulated “transmission” events with
1051 randomly assigned recipients. The black line represents the observed data, as in (D) the shaded
1052 regions represent the middle 50% and 95% of simulated outcomes. The results from the
1053 simulated logit models were smoothed by plotting the predicted probability of transmission at
1054 0.02 intervals. (F) The beta-binomial model fit. Similar to (D) except the simulated outcomes are
1055 the based on a beta-binomial model using a mean bottleneck of 1.71.

1056

1057 **Figure 4.** Combined estimates of within host mutation rate and effective population size.

1058 Contour plot shows the log likelihood surface for estimates of the effective population size and
1059 neutral mutation rate. The point represents the peak ($\mu = 4 \times 10^{-6}$, $N_e = 33$, log likelihood = -

1060 3,271). Log likelihoods for each contour are indicated.

1061

1062 **Supplemental Figure 1.** Sequence coverage for all samples. For each sample, the sliding
1063 window mean coverage was calculated using a window size of 200 and a step of 100. The
1064 distributions of these means are plotted as box plots (median, 25th and 75th percentiles,
1065 whiskers extend to most extreme point within median \pm 1.5 x IQR) where the y-axis represents
1066 the read depth and the x-axis indicates the position of the window in a concatenated IAV
1067 genome.

1068

1069 **Supplemental Figure 2.** Approximate maximum likelihood trees of the concatenated coding
1070 sequences for high quality H1N1 samples. The branches are colored by season; the tip
1071 identifiers are colored by household. Arrows with numbers indicate consensus and putative
1072 minor haplotypes for samples with greater than 10 iSNV.

1073

1074 **Supplemental Figure 3.** Approximate maximum likelihood trees of the concatenated coding
1075 sequences for high quality H3N2 samples. The branches are colored by season; the tip
1076 identifiers are colored by household. Arrows with numbers indicate consensus and putative
1077 minor haplotypes for samples with greater than 10 iSNV.

1078

1079 **Supplemental Figure 4.** The effect of titer and vaccination on the number of iSNV identified.
1080 (A) The number of iSNV identified in an isolate (y-axis) plotted against the titer (x-axis,
1081 genomes/ μ l transport media). (B) The number of iSNV identified in each isolate stratified by
1082 whether that individual was vaccinated or not. Red bars indicate the median of each distribution.

1083

1084 **Supplemental Figure 5.** Minority nonsynonymous iSNV in global circulation.

1085 The global frequencies of the amino acids that were found as minority variants in sample

1086 isolates (x-axis) plotted overtime (y-axis). Each amino acid trace is labeled according to the H3
1087 number scheme. All samples were isolated in December of 2014 (gray line).

1088

1089 **Supplemental Figure 6.** Reproducibility of iSNV identification for paired samples acquired on
1090 the same day. The x-axis represents iSNV frequencies found in the home-acquired nasal swab.

1091 The y-axis represents iSNV frequencies found the clinic-acquired combined throat and nasal

1092 swab.

1093

1094 **Table 1. Influenza viruses over five seasons in a household cohort**

	2010-2011	2011-2012	2012-2013	2013-2014	2014-2015
Households	328	213	321	232	340
Participants	1441	943	1426	1049	1431
Vaccinated, n (%) ^a	934 (65)	554 (59)	942 (66)	722 (69)	992 (69)
IAV Positive Individuals ^b	86	23	69	48	166
H1N1	26	1	3	47	0
H3N2	58	22	66	1	166
IAV Positive Households ^c					
Two individuals	13	2	9	7	23
Three individuals	5	2	3	3	11
Four individuals	-	-	1	2	4
High Quality NGS Pairs ^d	4	1	2	6	39

1095 ^a Self reported or confirmed receipt of vaccine prior to the specified season.

1096 ^b RT-PCR confirmed infection.

1097 ^c Households in which two individuals were positive within 7 days of each other. In cases of trios and quartets, the putative chains could have no pair with onset >7 days apart.

1098 ^d Samples with >10³ genome copies per µl of transport medium, adequate amplification of all 8 genomic segments, and average sequencing coverage >10³ per nucleotide.

1095
1096
1097
1098
1099
1100
1101
1102

1103 **Table 2. Within host effective population size of IAV**

Model	SNV Used	Generation Time (h)	Effective Population Size (95% CI)
Diffusion approximation	All	6	35 (26-46)
	All	12	17 (13-23)
Discrete model	All	6	32 (28-41)
	Nonsynonymous	6	30 (21-40)
	Synonymous	6	37 (27-54)
	All	12	23 (23-29)
	Nonsynonymous	12	19 (19-21)
	Synonymous	12	27 (22-33)

1104

1105

1106

1107 **Supplemental Table 1. Sensitivity and specificity of variant detection**

1108

Copy Number ^a	Variant Frequency	Original Pipeline ^b		Current Pipeline ^c	
		Sensitivity	Specificity	Sensitivity	Specificity
>10 ⁵	0.05	1	>0.9999	0.85	1.000
	0.02	0.85	0.9999	0.15	1.000
	0.01	0.95	0.9995	-	-
	0.005	0.35	0.9999	-	-
10 ⁴ -10 ⁵	0.05	0.95	0.9999	0.85	1.000
	0.02	0.9	0.9999	0.15	1.000
	0.01	0.8	0.9998	-	-
	0.005	0.4	0.9999	-	-
10 ³ -10 ⁴	0.05	0.8	>0.9999	0.70	1.000
	0.02	0.45	0.9999	0.15	1.000
	0.01	0.2	0.9997	-	-
	0.005	0.1	0.9999	-	-

1109

1110

1111

1112

^a Per μ l transport media

^b As described in McCrone JT and Luring AS, J. Virol. 90(15):6884, 2016.

^c As described in Methods, benchmarked for frequencies 0.02-0.98 only

1113
1114
1115

Supplemental Table 2. Nonsynonymous substitutions in HA antigenic sites

House ID	Enrolment ID	Symptom Onset	Subtype	Frequency	Amino Acid Change	Antigenic Site	Vaccinated	Day of Symptoms
1111	300481	3-30-2011	H3N2	0.071	E62G	E*	No	0
2166	320661	2-13-2012	H3N2	0.071	V297A	C	Yes	1
1302	301355	3-20-2011	H3N2	0.088	L86I	E	Yes	1
3075	331045	12-10-2012	H3N2	0.066	I214T	D	Yes	1
5219	50935	12-5-2014	H3N2	0.175	F193S	B*†	No	3
5263	51106	12-6-2014	H3N2	0.111	T128A	B	Yes	3
5290	51225	12-15-2014	H3N2	0.405	I260V	E*	Yes	1
5302	51273	12-13-2014	H3N2	0.030	S262N	E*	Yes	0
5098	50419	12-22-2014	H3N2	0.364	G208R	D	Yes	4
5033	50141	12-3-2014	H3N2	0.032	A163T	B	Yes	2
5034	50143	1-11-2015	H3N2	0.119	I307R	C	Yes	1
5289	51220	12-13-2014	H3N2	0.038	K189N	B*†	Yes	-1
5033	50141	12-3-2014	H3N2	0.025	D53E	C*	Yes	1
5033	50141	12-3-2014	H3N2	0.023	S312G	C	Yes	1
5269	51132	12-6-2014	H3N2	0.028	I242T	D	Yes	2
5147	50630	11-18-2014	H3N2	0.164	I242L	D	Yes	1
5034	50143	1-11-2015	H3N2	0.161	I307R	C	Yes	2
4185	UM40738	12-14-2013	H1N1	0.021	R208K	Ca	No	2

1116
1117
1118

* Sites observed to vary between antigenically distinct strains in Wiley et al., 1981 and Smith DJ et al., 2004.

† Sites located in the “antigenic ridge” identified in Koel et al., 2013.

Figure 1

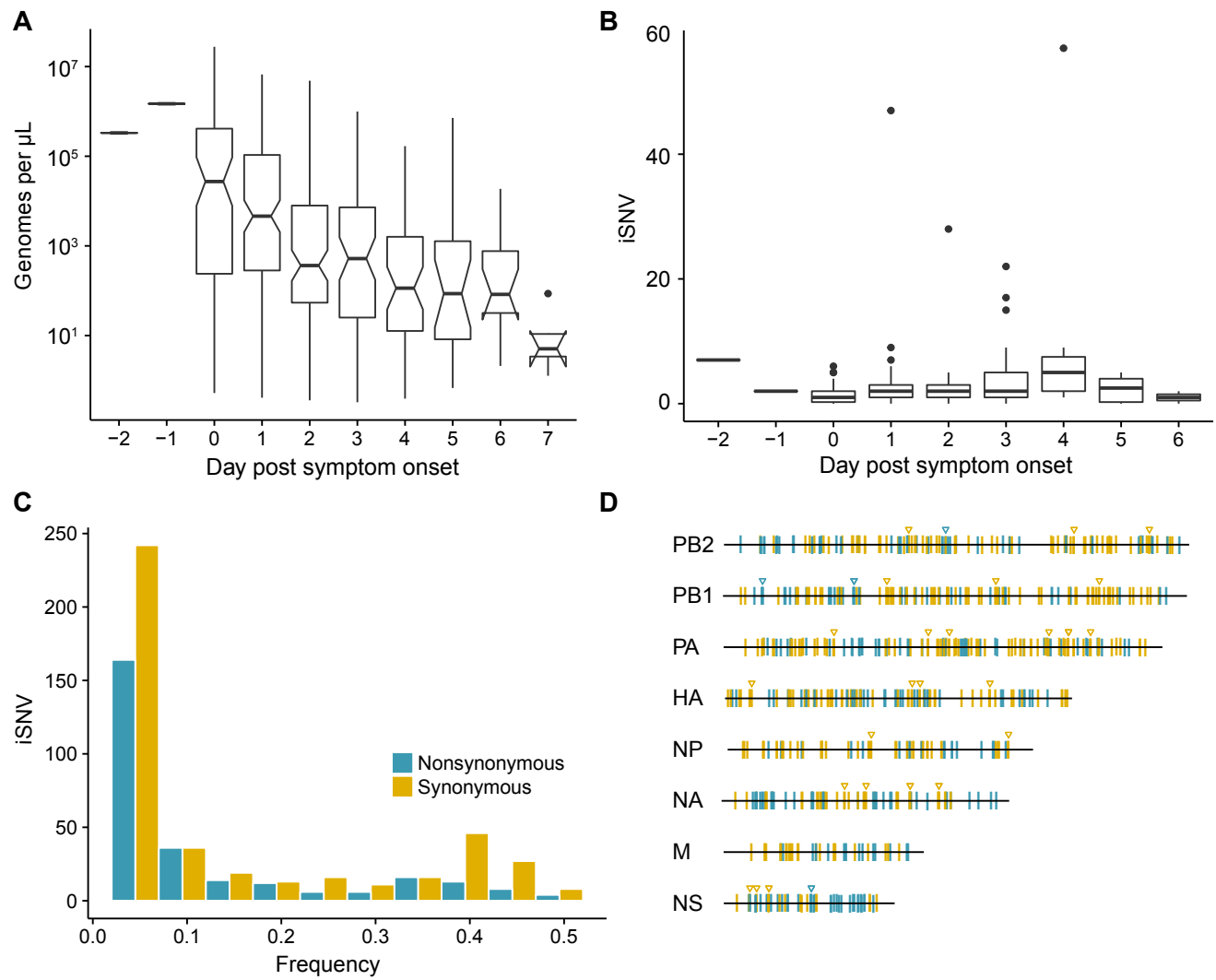


Figure 2

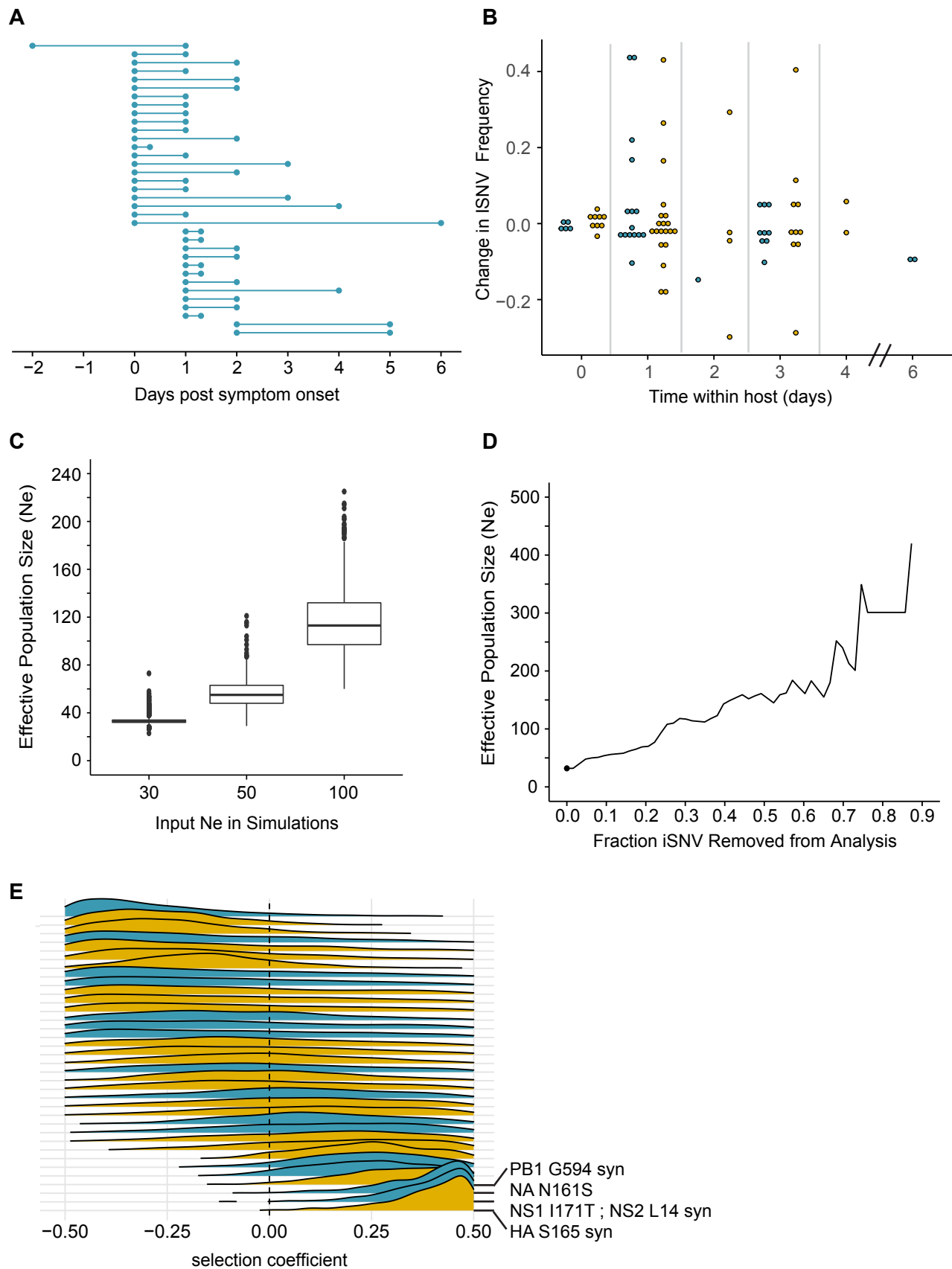


Figure 3

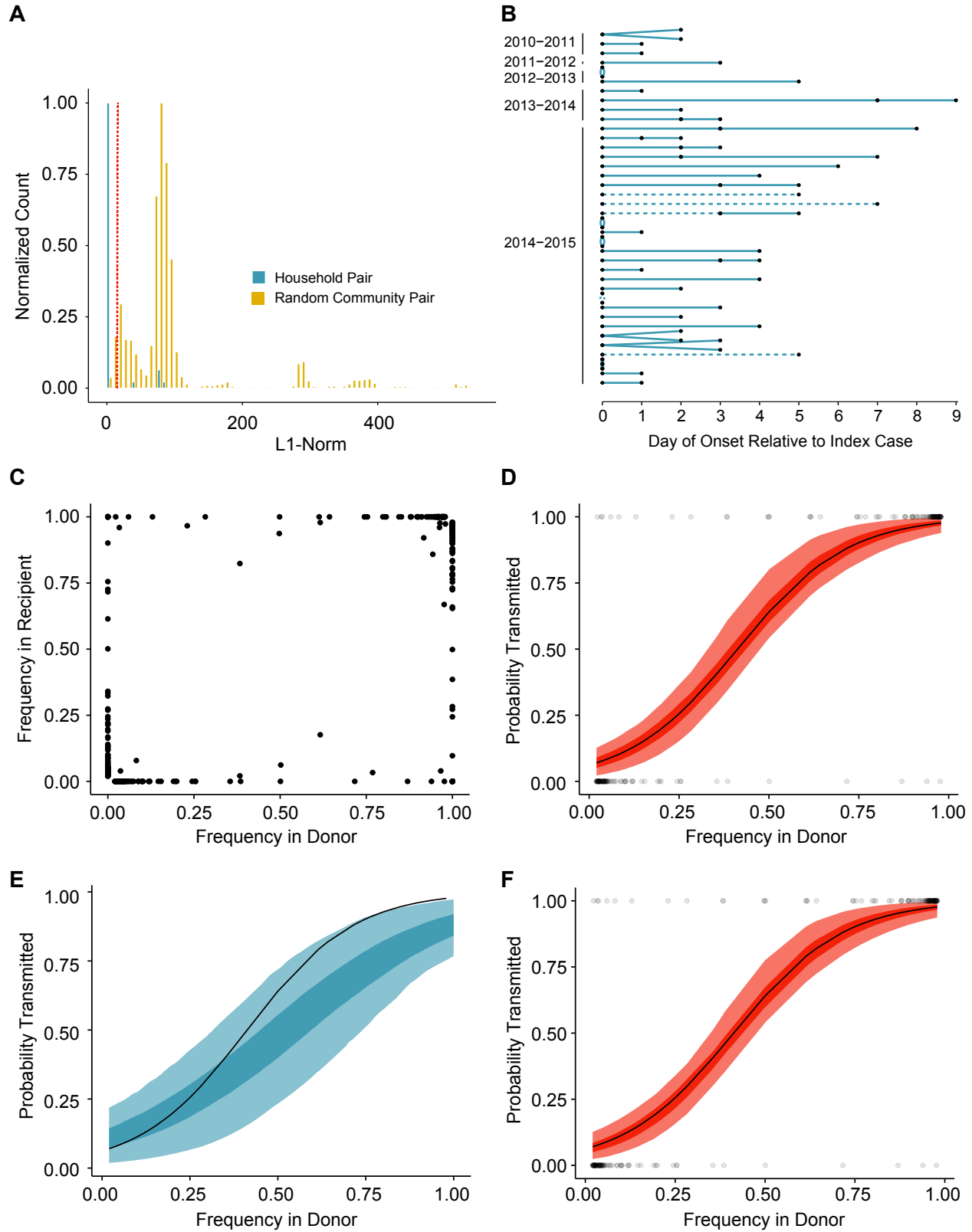


Figure 4

