

# 1 Timing HIV infection with nonlinear viral 2 dynamics

3 Daniel Reeves\*[1], Morgane Rolland [2,3], Bethany L Dearlove [2,3], Yifan Li [2,3], Merlin Robb [2,3],  
4 Joshua T Schiffer [1,4], Peter Gilbert [1,5], E Fabian Cardozo-Ojeda<sup>†</sup> [1], Bryan Mayer<sup>†</sup> [1]

5  
6 [1] Vaccine and Infectious Diseases Division, Fred Hutchinson Cancer Research Center, Seattle WA, USA  
7 [2] U.S. Military HIV Research Program, Walter Reed Army Institute of Research, Silver Spring, MD, USA.  
8 [3] Henry M. Jackson Foundation for the Advancement of Military Medicine, Bethesda, MD, USA.,  
9 [4] Department of Medicine, University of Washington, Seattle, WA, USA.  
10 [5] Department of Statistics, University of Washington, Seattle, WA, USA.

11 \*Corresponding author [dreeves@fredhutch.org](mailto:dreeves@fredhutch.org)

12 <sup>†</sup>These authors contributed equally

## 13 Abstract

14 In HIV prevention trials, precise identification of infection time is critical to quantify drug efficacy but  
15 difficult to estimate as trials may have relatively sparse visit schedules. The last negative visit does not  
16 guarantee a boundary on infection time because viral nucleic acid is not present in the blood during  
17 early infection. Here, we developed a framework that combines stochastic and deterministic within-host  
18 mathematical modeling of viral dynamics accounting for the early unobservable viral load phase until it  
19 reaches a high chronic set point. The infection time estimation is based on a population non-linear  
20 mixed effects (pNLME) framework that includes the with-in host modeling. We applied this framework  
21 to viral load data from the RV217 trial and found a parsimonious model capable of recapitulating the  
22 viral loads. When adding the stochastic and deterministic portion of the best model, the estimated  
23 infection time for the RV217 data had an average of 2 weeks between infecting exposure and first  
24 positive. We assessed the sensitivity of the infection time estimation by conducting *in silico* studies with  
25 varying viral load sampling schemes before and after infection. pNLME accurately estimates infection  
26 times for a daily sampling scheme and is fairly robust to sparser schemes. For a monthly sampling  
27 scheme before and after first positive bias increases to -7 days. For pragmatic trial design, we found  
28 sampling weekly before and monthly after first positive allows accurate pNLME estimation. Our  
29 estimates can be used in parallel with other approaches that rely on viral sequencing, and because the  
30 model is mechanistic, it is primed for future application to infection timing for specific interventions.

## 31 Introduction

32 A key challenge for HIV prevention trials is to identify the timing of the exposure that ultimately led to  
33 breakthrough infection. Estimation of infection time subsequently allows inference of the concentration  
34 of the protective agent at exposure, which is critical to understanding why HIV acquisition was not  
35 prevented. Early infection is difficult to study in practice; even if prospective sampling were available,  
36 HIV RNA is not detectable in blood during early HIV infection and participants cannot necessarily point  
37 to potential recent exposure events with accuracy. Therefore, to estimate time of infection, a model or  
38 inference technique is required.

39 Estimation techniques have been described previously. Several use viral sequence data and evolutionary  
40 models to trace time back to the founder sequence<sup>1-3</sup>. Others use viral load data prior to viral peak and  
41 retrace using log-linear regression (average or maximum upslope)<sup>3</sup>. Others apply diagnostic ‘window  
42 times’ that leverage Fiebig staging<sup>4</sup> and prior knowledge of eclipse phase duration<sup>5,6</sup>, where eclipse  
43 phase is defined as the period of time between HIV acquisition and first detectable viral load. Finally,  
44 combinations of these approaches have been organized into a statistical framework<sup>7</sup>.

45 Here, we introduce an approach applying such viral dynamics models with statistical inference on viral  
46 load data. Model development was achieved through fitting to longitudinally sampled viral loads from  
47 46 participants in the RV217 study<sup>8</sup>. Population nonlinear mixed-effects (pNLME) modeling was used to  
48 determine the optimal model parameters for each individual, given a population distribution. We tested  
49 30 candidate models informed by previous viral load modeling, and the best model was selected by  
50 parsimony. The very early moments of HIV infection are thought to be stochastic<sup>9</sup>, and have been  
51 modeled with stochastic viral dynamics<sup>10,11</sup>. Therefore, we used our best model in a stochastic  
52 formulation to simulate early HIV dynamics, allowing for fluctuations and extinction by chance.  
53 Together, the stochastic and deterministic models provide an estimate and associated uncertainty  
54 interval for the infection time of each individual in the study.

55 The pNLME modeling approach provides several advantages. By using a population model, it is possible  
56 to estimate infection times in individuals with sparse viral load data, including those without any  
57 measurements during viral upslope. Viral dynamics are not as sensitive to multiple founder infections as  
58 evolutionary methods. And finally, mechanistic models have been used to describe viral dynamics during  
59 broadly neutralizing antibodies therapies<sup>12</sup>, suggesting our methodology might be applicable to infection  
60 time estimation from emerging trial data including a therapeutic prevention modality.

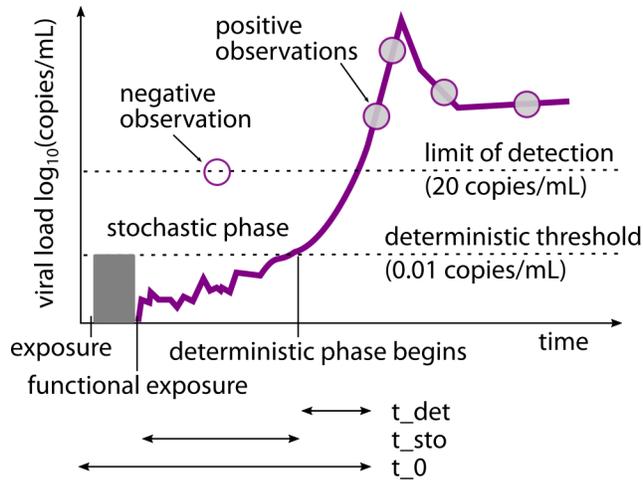
## 61 Results

62 **A framework for estimating infection time using viral dynamics.** Using experimental data and  
63 modeling, we set out to develop a framework for estimating HIV infection time from viral load data.  
64 Using observed first positive viral load, we worked backwards toward infection time and defined several  
65 precise moments during HIV primary infection for modeling (**Fig 1**). HIV infection begins with an  
66 *infecting exposure*, the target time of estimation. Starting with this exposure event, there is a brief  
67 “black-box” phase encompassing biology not captured with past viral dynamic models. For example, the  
68 virus may need to diffuse or clear mucosal and anatomical barriers before beginning viral replication as  
69 described by mechanistic models. Animal challenge studies and human cases where infecting exposure  
70 is almost certainly known suggest this period is brief, from a few hours to 1 day<sup>9,13,14</sup>, but given the lack  
71 of information we note it here as a fundamental uncertainty and potential bias in our estimates.

72 Next, we assumed that viral kinetics can be described by a dynamic model unifying observable and  
73 unobservable viral loads. At this point, we assumed bottlenecking has resulted in at most a few infected  
74 cells starting a productive infection in the human body. Viral expansion from these few cells begins a  
75 stochastic phase lasting a duration defined as the stochastic phase,  $t_{sto}$ . The stochastic phase likely  
76 encompasses the initial viral replication at the infection foci and the transition from a localized infection  
77 to an infection that has reached germinal centers. The stochastic phase ends when the viral load crosses  
78 a *deterministic threshold*. We estimated this threshold through repeated stochastic simulations finding  
79 the minimum viral load where the 1) the slope of stochastic viral loads were nearly log-linear and 2)  
80 there was effectively no chance of stochastic burn out. We determined that a value of 0.01 copy/mL  
81 sufficiently satisfied both criteria. The time between viral load crossing the deterministic threshold and  
82 reaching the first positive viral load observation was defined as the *deterministic phase*  $t_{det}$ . Finally, the

83 time between infecting exposure and first positive viral load, comprising these three phases, we  
84 cumulatively refer to as  $t_0$ . This time interval has been referred to as the *eclipse phase*<sup>5</sup>.

85

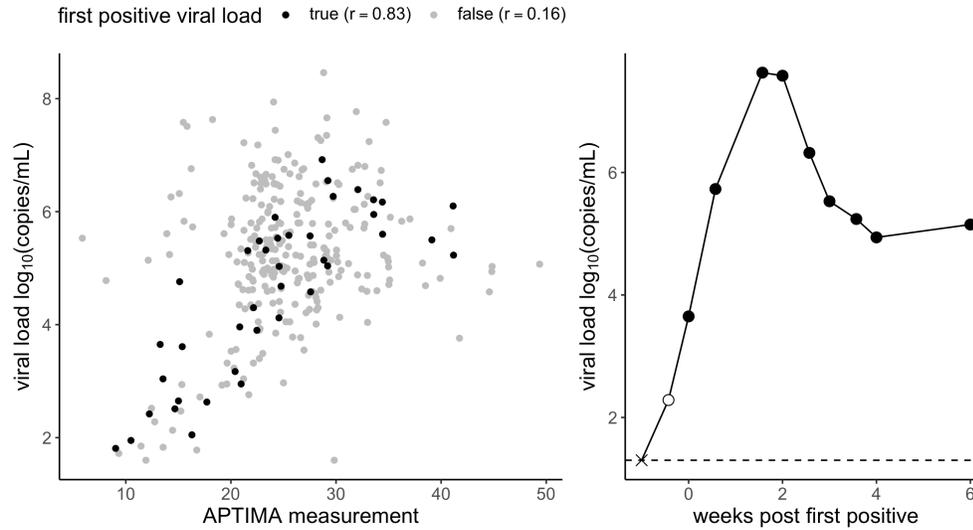


86

87 **Fig 1. Cartoon schematic of modeling definitions.** The time between infecting exposure and first positive  
88 viral load can be described in 3 phases. First, we recognize the possibility of an unknown but likely brief  
89 “black-box” period describing localized biology that exists immediately following infecting exposure. We  
90 assume this period is short compared to the following phases. Second, a stochastic process governs early  
91 viral expansion, starting with one or a few infected cells initiating systemic infection in the new host –  
92 and concluding when viral load reaches the deterministic expansion threshold ( $t_{sto}$ ). Third, a  
93 deterministic model ( $t_{det}$ ) proceeds, describing the observed viral dynamics. By combining estimates for  
94 these phases, we finalize our estimate of  $t_0$ , the time between infecting exposure and first positive viral  
95 load, sometimes referred to as the *eclipse period*.

96 **Experimental data for model development.** We used viral load observations from the RV217 study<sup>8</sup>  
97 including 46 individuals out of 155 total diagnosed acute HIV-1 infections in the study. Individuals had  
98 twice-weekly HIV tests before diagnosis using the APTIMA HIV-1 RNA Qualitative Assay (Hologic)—a  
99 fingerstick device testing small blood collection (0.5 mL). Once diagnosed (2 APTIMA positive visits),  
100 quantitative PCR was used to quantitate HIV RNA twice weekly in these individuals, who did not initiate  
101 antiretroviral treatment (ART) and had ~10 study visits in the first month after diagnosis. From this  
102 cohort, we assembled viral loads from Thai and Ugandan men, women and transgender individuals. Only  
103 individuals with more than 3 detectable longitudinal viral load observations were included. We found  
104 that in very early infection, APTIMA and viral load were strongly correlated (**Fig 2A**) and APTIMA  
105 measurements could be used to impute viral load at diagnosis times for individuals without measured  
106 viral loads (**Fig 2B**). Using this relationship, we imputed viral load at APTIMA diagnosis for 28  
107 participants, which adjusted the first positive viral load by a few days. Several individuals were not  
108 diagnosed until later acute infection, meaning that peak and upslope of viral load are not obviously  
109 detected. We do not exclude these individuals, instead relying on our population modeling approach  
110 and borrowing strength across the cohort to make estimates. These estimates are particularly useful  
111 because such data sets provide significant challenges with other modalities.

112



113

114 **Fig 2. Correlation plot between APTIMA measurements and viral load.** A) A strong linear correlation  
115 (Pearson  $r = 0.83$ ) is found between APTIMA and log viral load at the first positive viral load (black dots).  
116 If positive samples beyond the first positive (gray dots) are included, and at higher measurements of  
117 either outcome, APTIMA is less correlated to viral load. B) We used diagnostic APTIMA measurements  
118 prior to first positive viral load to impute additional viral load values. Here the closed circles indicate  
119 observed viral load measurements, the 'x' the last negative measurement, and the open circle indicates  
120 an APTIMA imputed value.

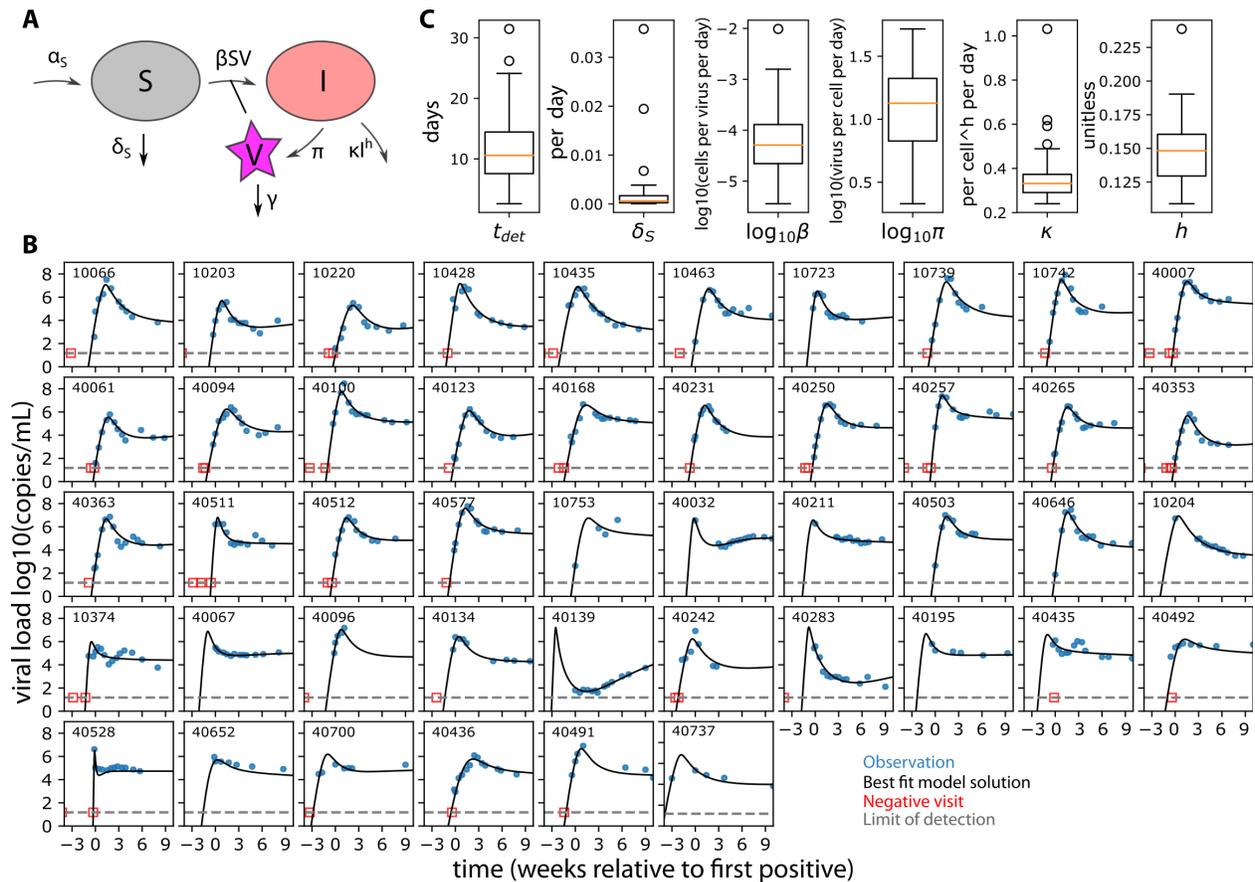
121 **Inference of  $t_{det}$  from a parsimonious model to the RV217 cohort data.** The first step in estimating  $t_{det}$   
122 was developing a model that best-described the observed data. Thus, we selected four distinct and  
123 previously applied mechanistic models of HIV primary infection and varied their population  
124 parameterizations (the number and type of parameters estimated). This resulted in a total of 30 models  
125 (see **Supplementary Table 1**). The four mechanistic models included the canonical viral dynamics  
126 model<sup>15</sup>, two models recently fit to SHIV/SIV viral dynamics<sup>16,17</sup>, and our own simplified model based  
127 upon Ref<sup>18</sup>. We found that the most parsimonious model to the RV217 cohort data (**Supplementary**  
128 **figure 1 and Table 1**) includes susceptible target cells ( $S$ ) that are born and die naturally and virus ( $V$ )  
129 that infects these cells and creates productively infected cells that produce viable virus ( $I$ ). Infected cell  
130 death rate depends on their own density powered by an exponent ( $h$ ). This term semi-mechanistically  
131 encapsulates natural cytopathic cell death during viral production, as well as innate or acquired  
132 immunity against HIV infected cells that escalates as the number of infected cells increases (**Fig 3A**, see  
133 also **Methods Eq. 1**). In this way, an explicit immune effector compartment is not needed, and the  
134 model is simplified substantially.

135 The model output is congruent with previous data for other model compartments. For example, it  
136 predicts a susceptible cell drop between 40—80%<sup>19</sup> (which may relate to the CD4+ T cell depletion  
137 during peak viremia<sup>20</sup>) and allows for the large observed inter-participant variation of viral peak  
138 (**Supplementary figure 2A**).

139 The best fit model for each individual is displayed in **Fig 3B**. We used population nonlinear mixed-effects  
140 (pNLME) modeling to estimate parameters, such that each individual has their own estimated  
141 parameters, but these estimates are constrained to be drawn from population distributions of each  
142 parameter; the population distribution is simultaneously estimated. All distributions of parameter  
143 estimates are shown in **Fig 3C** and values are quoted for each individual in **Supplementary Table 2**.

144 From this model, the estimated time between deterministic threshold and first positive viral load,  $t_{det}$ ,  
 145 (Fig 1) ranged from 2.5–32.6 days across the 46 participants with a median of 10.1 days. We also  
 146 verified that models with comparable AIC (<10 difference from the best model AIC score) predict similar  
 147 individual values for  $t_{det}$ . Summary statistics of viral load (peak and set point) were not correlated with  
 148 deterministic time  $t_{det}$ ; rather they were strongly correlated with estimated infectivity ( $\beta$ ), viral  
 149 production rate ( $\pi$ ) and the nonlinear death exponent ( $h$ ) (Supplementary figure 2B). The magnitude of  
 150 the first positive viral load was significantly, but not strongly, correlated with  $t_{det}$  (Supplementary figure  
 151 3). These results suggest other estimated parameters are mostly independent of infection timing and  
 152 that the model predictions are informative beyond upslope regression—i.e. nonlinear estimation  
 153 enhances our predictive power.

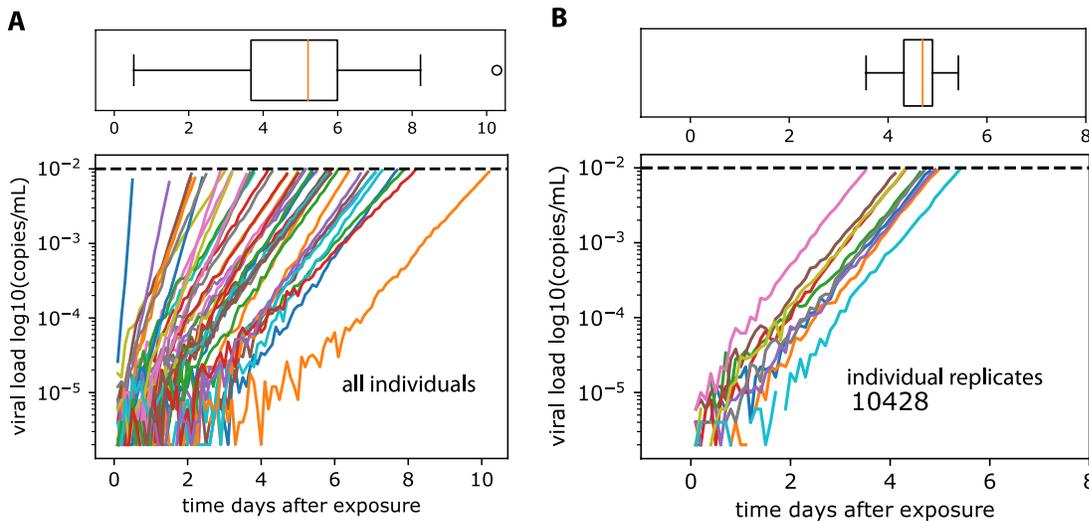
154



155

156 **Fig 3. The optimal mathematical model recapitulating RV217 viral load kinetics.** A) By comparison of  
 157 four structurally distinct models and many distinct statistical population models for each, we effectively  
 158 tested 30 models for the data and arrived at an optimally parsimonious model, schematized here. The  
 159 model is identical to the canonical viral dynamics model except infected cells have a nonlinear death rate  
 160 (see Eq 1). B) This model recapitulates diverse viral load kinetics in the RV217 human study. In each  
 161 panel, viral data are gray dots and best individual fit is a blue line. By borrowing strength through the  
 162 population fitting approach, the model infers peak and upslope even when those data are missing (see  
 163 40139, or 40700 for example). Last negative visits are shown as viral loads at the limit of detection (20  
 164 copies/mL) and were included as censored data for fitting. Only 1 individual (last panel, 40737) had a  
 165 first positive viral load that would be shifted substantially given the APTIMA imputation. C) 6 parameters  
 166 were estimated including the deterministic infection time, the crucial variable for timing infection.

167 **Stochastic simulations until the deterministic threshold.** Evidence from modeling other viruses suggests  
168 that early stochastic events are linked to later deterministic kinetics<sup>21</sup>. For example, for cytomegalovirus  
169 (CMV) infection, extinction probabilities, duration, and magnitude of transient stochastic infections are  
170 consistent with primary infection mathematical model parameters<sup>22</sup>. Therefore, based on the individual  
171 best fit parameter sets, we performed stochastic simulations to determine the time-window between  
172 the introduction of a single infected cell and the deterministic threshold ( $t_{sto}$  in **Fig 1**). Simulations were  
173 initialized with a single infected cell per  $\mu\text{L}$   $I(0) = 1$  and at the viral free equilibrium between  
174 susceptible cell birth and death  $S(0) = vol \times \alpha_S / \delta_S$ . Scaling up to realistic volumes allows for a  
175 discretized stochastic simulation;  $vol$  was chosen to be  $5 \times 10^8 \mu\text{L}$ , or 5 L of blood (typical for adult  
176 human) at approximately 100-fold concentration based on the finding that the majority of lymphocytes  
177 reside in lymphoid tissues where infection is assumed to initiate before spilling over into blood<sup>9,23</sup>.



178

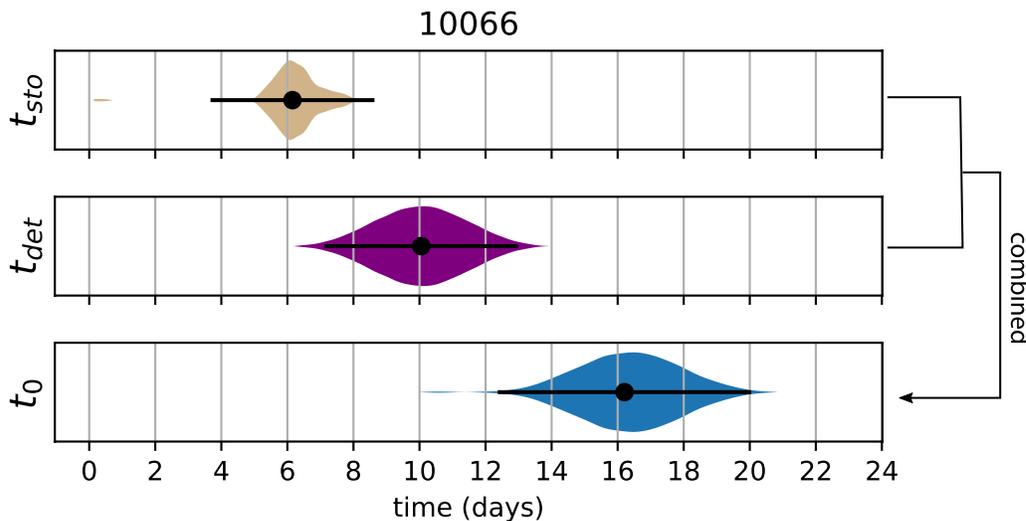
179 **Fig 4. Stochastic simulations using best-fit parameter estimates from deterministic model.** Viral load  
180 kinetics until deterministic threshold (0.01 copy/mL). A) A single stochastic realization for all 46 individual  
181 parameter sets from **Fig 3** with associated boxplot of the distribution of times (of these 46 simulations)  
182 to reach the deterministic threshold. The slopes are different across individuals owing to the different  
183 parameter estimates from the deterministic models. B) 10 replicate stochastic realizations for a single  
184 individual until deterministic threshold with associated boxplot of the distribution of times (of these 10  
185 simulations) to reach the deterministic threshold. Here slopes are nearly identical, but due to the  
186 stochasticity of the simulation, the time to reach the deterministic threshold varies between 3-5 days.  
187 Note discontinuities in lines are artifacts of downsampling for file size considerations.

188 For each individual, the best fit parameters of the deterministic model were used to conduct 10  
189 stochastic simulations via the tau-leap method<sup>24</sup>. Because HIV transmission is a rare per coital event<sup>25</sup>  
190 and we are interested in infection time estimation, we conditioned upon successful infection<sup>10</sup> by only  
191 using simulations from stochastic runs that did not go extinct. The simulations were halted when viral  
192 load crossed the deterministic threshold (0.01 copies/mL) and the time to reach that viral level ( $t_{sto}$ )  
193 was recorded. Simulated viral loads from a single stochastic simulation of each individual are shown in  
194 **Fig 4A**. The distribution of stochastic times ( $t_{sto}$ ) is visualized above the viral load panel, indicating a  
195 slightly asymmetric time to crossing the deterministic threshold with median  $\sim 5$  days in this single  
196 stochastic simulation. There is substantial variability in the slope of these viral load trajectories based on  
197 the range of parameters inferred from the deterministic model for each individual. We also performed  
198 replicate simulations for single individuals (10 replicates for participant 10428 are shown in **Fig 4B**). In

199 this case, viral load slopes are nearly identical by the time the deterministic threshold is crossed, but the  
200 early stochastic events introduce some variability in  $t_{sto}$ . For this individual, the median time between  
201 infection and deterministic threshold was 5 days, with total range between 3-5 days in these 10  
202 simulations. In summary, viral load upslope varies highly across subjects but minimally within-subjects.  
203 Variability introduced by the stochastic phase is predominantly a shift, rather than a scaling of infection  
204 time. This agrees with modeling of barcoded virus data early in infection (recently reported by Docken et  
205 al. during Dynamics & Evolution of HIV and Other Viruses 2020).

206 An important parameter for these simulations is the initial number of infected cells. We show estimates  
207 of  $t_{sto}$  are inversely correlated with  $I(0)$ . For example, as  $I(0)$  was increased from 1, 10, 100, to 1000,  
208 the median estimate of  $t_{sto}$  across individuals decreased 5-1 days (**Supplementary figure 4**). As a result,  
209 this difficult-to-measure biological parameter only adjusts estimates by a few days.

210 **Combining the stochastic and deterministic phases to estimate infection time.** Next, we integrated the  
211 stochastic and deterministic timing estimates to complete the estimation of  $t_0$ , the time between  
212 infecting exposure and first positive viral load, or the eclipse phase (**Fig 1**). Here as an example, we  
213 present this procedure for individual 10066 (**Fig 5**). First, we used the best fit parameters and performed  
214 100 replicate stochastic simulations to estimate a distribution of  $t_{sto}$ ; the mean was approximately 6  
215 days, and the distribution was skewed, with 95% uncertainty interval ranging between 4-9 days. Second,  
216 we drew values of  $t_{det}$  from a constructed conditional distribution using Markov-Chain Monte-Carlo  
217 given the population and random effect estimates of  $t_{det}$  (mean 10, 95% uncertainty interval between  
218 7-13 days). The infection time,  $t_0$ , was then calculated by drawing and summing  $t_{sto}$  and  $t_{det}$  from their  
219 respective distributions. This was repeated 10000 times to generate an average  $t_0$  with associated 95%  
220 uncertainty interval (see estimates of  $t_{sto}$ ,  $t_{det}$ , and  $t_0$  for this individual in **Fig 5**). We estimated that  
221 this individual's infection occurred 16 days prior to first positive viral load with 95% uncertainty interval  
222 ranging between 12 and 20 days. This procedure was performed for all individuals.



223  
224 **Fig 5. Individual estimate example.** Bootstrap combination of the deterministic and stochastic  
225 estimates provides an estimate for an individual's (10066) time of infection. The stochastic time interval  
226 ( $t_{sto}$ ) between 1 infected cell and the deterministic threshold (0.01 copies/mL) was determined by 100  
227 replicate stochastic simulation for that individual. Here the mean estimate of  $t_{sto}$  was 6 days (dot) with  
228 95% uncertainty interval (lines) ranging between 4-9 days. The entire probability distribution is shown to  
229 illustrate skew. The time interval between the deterministic threshold and the first positive viral load was

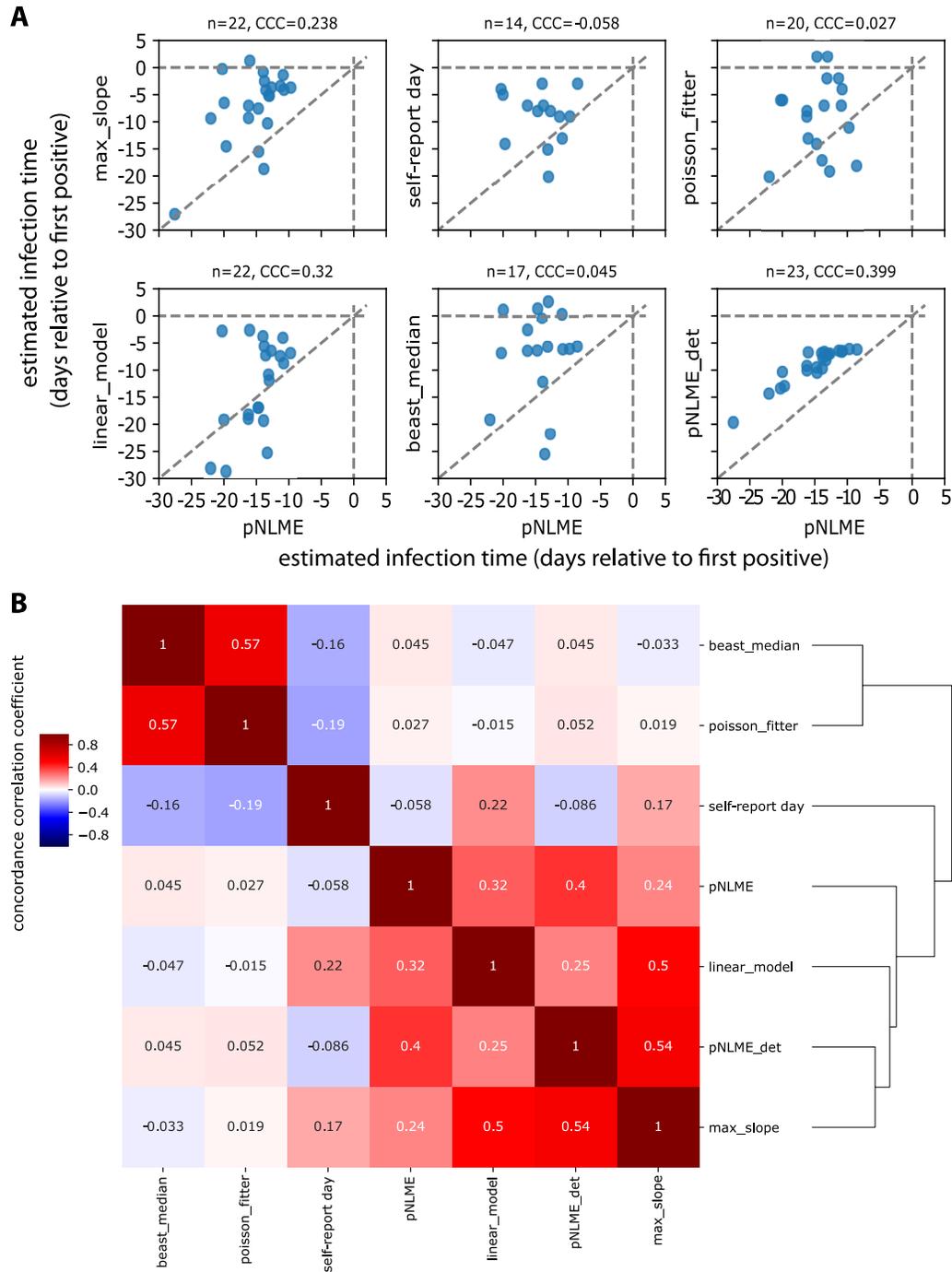
230 *determined by the best estimate of  $t_{det}$  using population nonlinear mixed effects modeling (see Fig 3).*  
231 *Here the mean (dot, ~10 days) and 95% uncertainty interval (lines, ranging between 7-13 days) from the*  
232 *MCMC estimate are shown with the derived distribution. Finally, the time between infection and first*  
233 *positive viral load ( $t_0$ ) is calculated by 10000 random combinations of  $t_{det}$  and  $t_{sto}$  drawn from these*  
234 *distributions. Our best estimate suggests this individual was infected 16 days prior to first positive with*  
235 *95% uncertainty interval ranging from 12-20 days.*

236 **Direct comparison to previously applied infection timing estimation tools.** Rolland et al. used several  
237 viral load and phylogenetic inference techniques to estimate infection times using the RV217 data<sup>3</sup>.  
238 These methods are the maximum slope of any two points on the upslope (max\_slope), the best log-  
239 linear regression slope (linear\_model), self-reported entries from trial participants (self\_report),  
240 Bayesian phylogenetic inference of median time to most-recent common ancestor (BEAST)<sup>26</sup>, and  
241 Poisson fitter<sup>27</sup> diversity estimator based on envelope sequences sampled at three time points in the  
242 first six months of infection. We compared our viral dynamics population non-linear mixed effects  
243 (pNLME) modeling approach estimations directly against all methods in Fig 6. In general, our  
244 deterministic estimates were in the same relative range of the other estimates. However, concordance  
245 correlation coefficients (CCC), which score how close data lie to the line  $y=x$ , are in general fairly weak  
246 ( $CCC < 0.4$ ) between pNLME and other methods (Fig 6A). This is driven by the fact that the complete  
247 estimator finds infection time earlier than most other methods, perhaps due to the additional stochastic  
248 phase. Adjusting the initial number of infected cells from 1 to 1000 (see Supplementary figure 4), or  
249 removing the stochastic phase decreases the average eclipse time, closer to previous estimates.  
250 However, we show correlation between pNLME with and without the stochastic component (final panel  
251 in Fig 6A) to illustrate this relationship is not necessarily linear.

252 We also compared all previous point estimates to one another (Fig 6B). No approaches were very  
253 strongly correlated by CCC. Hierarchical clustering of previous methods shows there are two distinct  
254 groups that include genetic estimators (BEAST, PFitter) and viral dynamic estimators (max-slope,  
255 linear\_model). Self-report diary entries and our method (pNLME) fall roughly in between.

256 **Wide applicability of pNLME to sparse data.** The pNLME approach is widely applicable to data that  
257 challenge other methods. It can provide estimates for individuals for whom viral upslope is completely  
258 undetected. It also does not produce large outliers and never estimated the time of infection to be after  
259 first positive, as Max-slope, BEAST, and PFitter do in a few cases. pNLME also does not appear to be  
260 sensitive to multiple founder infections (which are particularly difficult for genetic estimators). For  
261 example, Rolland et al. identified some individuals in this cohort infected with multiple founder viruses  
262 based on the sequence analysis; for these infections, estimates of time to most recent common ancestor  
263 often gave estimates preceding the date of last negative test by many months (reflecting divergence in  
264 the transmitting partner rather than divergence after transmission)<sup>3</sup>. Infection with multiple founders  
265 has been associated with higher set-point viral loads<sup>28</sup>. Therefore, we tested to see if our model  
266 parameter estimates were different in single versus multi-founder infections (as differentiated by  
267 Rolland et al., see Supplementary figure 5). We observed no obvious patterns distinguishing single and  
268 multi-founder participants and found no significant differences among our parameters (Mann-Whitney  
269  $p > 0.1$ ) but note the limited sample size with these data ( $n=9$  multiple founders in this set). While beyond  
270 the scope of this paper, applying multi-founder status as a pNLME model covariate might admit more  
271 power given the small sample size. Importantly, the estimate of the time of theoretically crossing the  
272 detection limit was not affected by the distinction of multiple founders, meaning our estimates are  
273 robust to this challenge for phylogenetic inference.

It is made available under a [CC-BY-ND 4.0 International license](https://creativecommons.org/licenses/by-nd/4.0/).



274

275 **Fig 6. Comparisons of the population non-linear mixed model estimation (pNLME) approach for**  
 276 **infection timing versus 5 other methods.** Methods include the maximum slope of any two points on the  
 277 upslope, the best log-linear regression slope (*linear\_model*), self-reported entries from trial participants,  
 278 Bayesian phylogenetic inference of median time to most-recent common ancestor (BEAST), Poisson fitter  
 279 diversity-based estimator, and our own approach using only the deterministic component. A) Best  
 280 estimate of each available individual from each method expressed as predicted infection time relative to  
 281 first positive (all estimates are tabulated in **Supplementary Table 2.**) An estimate above 0 (see dashed  
 282 lines), indicates unrealistic estimates of infection after first positive. While our method can be used on all  
 283 46 individuals, other approaches are constrained by features of the data (e.g. detection of upslope,

284 sequencing characteristics). Thus, comparison size  $n$  is denoted above each panel. Concordance  
285 correlation coefficients show individual agreement is generally weak:  $CCC = 1$  when all data lie on the  
286 line  $y = x$  (shown as a dashed line). B) Hierarchical clustering using concordance correlation coefficients  
287 indicate which methods give most similar estimates. Sequence based methods and viral dynamic  
288 methods fall into 2 main clusters, with self-report falling in the middle of these. pNLME agrees more  
289 strongly with other viral dynamic methods.

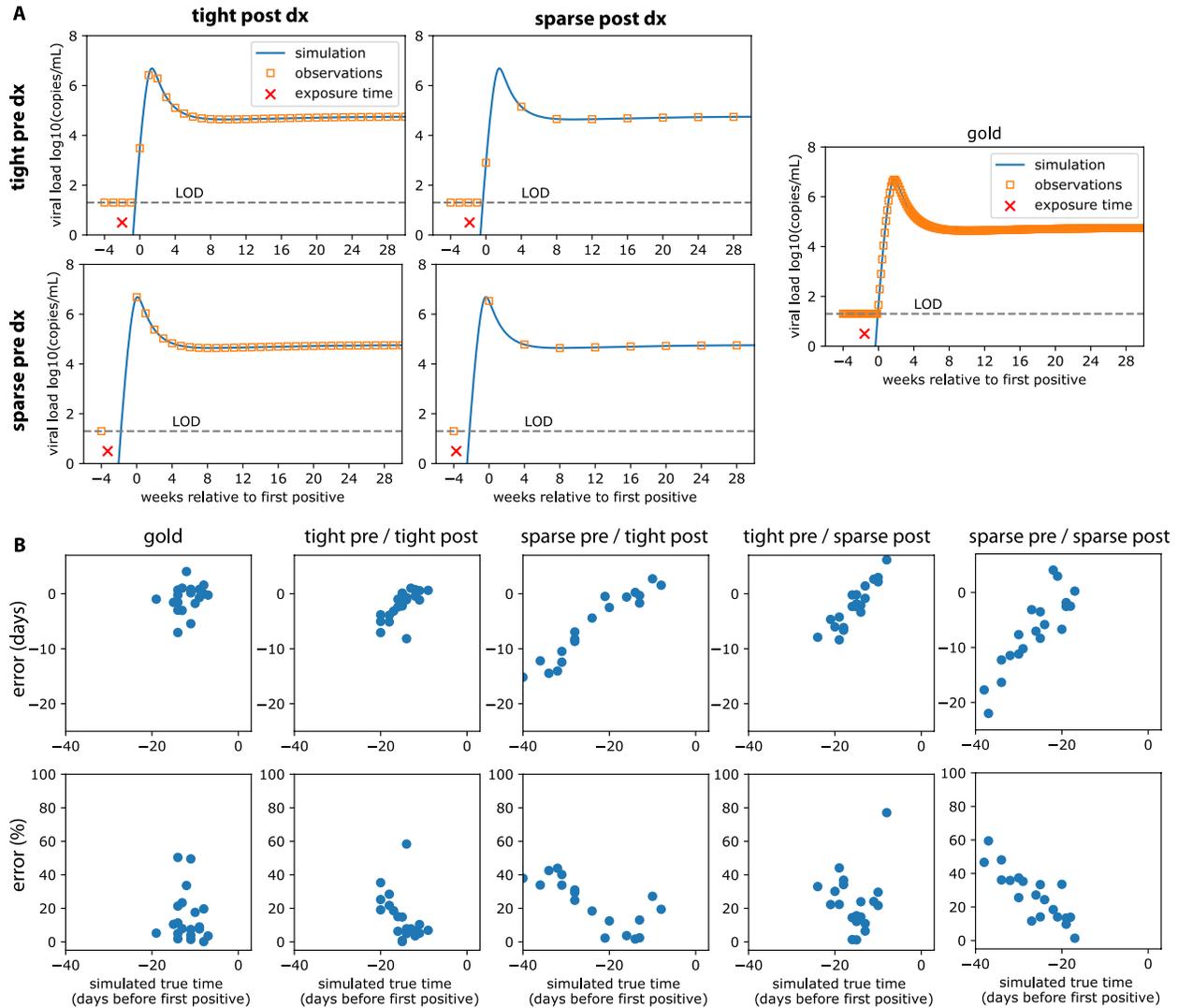
290 **Proof of concept study on synthetic data with realistic study protocols.** To assess the accuracy of  
291 pNLME estimates, we performed a simulation study using several different sampling schemes. We  
292 simulated viral loads from 20 randomly chosen RV217 participants and sampled these viral loads with 5  
293 different theoretical protocols. The first we refer to as “gold” meaning daily sampling before and after  
294 first positive. We refer to tight as weekly sampling visits, and sparse as monthly sampling visits (every 4  
295 weeks). Infection was assumed to occur uniformly between study visits. If viral load was above 20  
296 copies/mL at a visit that sample was called first positive (or diagnosis, dx), and measurements occurred  
297 subsequently. In **Fig 7A** we illustrate an example infection (red x), viral load (blue line) and observations  
298 (orange circles) for each protocol: “tight pre / tight post”, “tight pre / sparse post”, “sparse pre / tight  
299 post”, “sparse pre / sparse post”, and “gold”. We took these synthetic data observations and estimated  
300  $t_0$  with pNLME. In this step, we used the RV217-trained model, meaning that we fixed the population  
301 distributions (as we would with new test data), and arrived at a new conditional distribution of  
302 individual parameters for each synthetic data set. We applied those parameters to the stochastic  
303 modeling step, completing the estimate of  $t_0$  on each synthetic data set.

304 **Fig 7B** shows the absolute error (days difference between truth from the synthetic data and inferred  $t_0$   
305 from pNLME applied to those data) and the % error:  $(\text{true}-\text{inferred})/\text{true} \times 100\%$ . Gold standard and  
306 tight/tight predictably admitted the lowest errors. Very few individuals were overestimated, meaning  
307 that when inference was incorrect, their infection time was usually closer to first positive than inferred.  
308 The exception to this occurred for some individuals with sparse post sampling.

309 All schemes other than gold had an obvious bias. Absolute and percent error was higher in individuals  
310 for whom true infection time was earlier. This means that uncertainty rises with estimation time farther  
311 from first positive. Put another way, our confidence decreases as the estimator projects farther into the  
312 past— an intuitively satisfying, albeit challenging finding. That this effect was fairly linear hints that it  
313 might be corrected. However, this may be an artifact of our synthetic data exercise, so we opted not to  
314 follow through with any correction. Rather, we focus on individuals who appear to have been infected  
315 within 20 days since first positive. For all sampling schemes, error on these estimates has a median of  
316 +/-10 days. A corollary of this finding is that sparse sampling after diagnosis was less detrimental than  
317 sparse sampling before diagnosis, because of the growing uncertainty with time and the likelihood of  
318 missing upslope, peak, and downslope.

319 This exercise illustrates one of the most practical applications for this method: estimating infection time  
320 in clinical trials of HIV pre-exposure prophylaxis agents. The “sparse/sparse” case represents a protocol  
321 comparable to that of the AMP (antibody mediated prevention) studies. In that study, visits occur every  
322 4 weeks, and after first positive (week 0) visits occur at week 2, 4, 8, 12, and 24 weeks<sup>12,29,30</sup>. Thus, given  
323 the data generation distribution produced under our modeling assumptions, and the additional  
324 assumption that HIV dynamics in participants in the AMP study are comparable to participants in RV217,  
325 we expect our approach would provide reasonably accurate estimates for individuals who appear to  
326 have been infected within 20 days of first positive visit (95% uncertainty interval ~5 days), and less  
327 confident estimates for others. A secondary result of our modeling is that more sensitive detection  
328 could be crucial to avoiding the ultimately challenging case of individuals infected >4 weeks before first  
329 positive visit. APTIMA testing and viral load prediction as in **Fig 2** would help this substantially.

330



331

332 **Fig 7. Accuracy of timing tested on simulated viral load data with sparse and tight study sampling.** We  
 333 simulated viral loads using individual parameter sets and sampled according to 4 different theoretical  
 334 study protocols. We refer to tight as weekly, and sparse as monthly (every 4 weeks). We simulated data  
 335 assuming infection occurred uniformly throughout observations periods. If viral load was above the study  
 336 assay detection limit at an observation (20 copies/mL), that was called first positive (or diagnosis, dx),  
 337 and measurements occurred subsequently. A) 4 examples of each study protocol: combinations of tight  
 338 and sparse, pre and post diagnosis (dx). True infection is denoted with the red x, simulated viral load with  
 339 the blue line, and observations given the protocol with the orange squares. B) 20 individuals were  
 340 simulated with each protocol, and pNLME inference was performed on those data. The accuracy of the  
 341 estimated  $t_0$  compared to the true infection time is shown as error in days (difference between pNLME  
 342 estimated and true time) and percent error (error relative to true time x 100%) for each sampling  
 343 protocol. Low error therefore indicates estimates that agree, and percent error illustrates how estimates  
 344 get more biased (relatively to other estimates) as the time between detection and true infection time  
 345 increases.

## 346 Discussion

347

348 Estimating infection time is especially critical in HIV prevention trials. If drug levels at the precise time of  
349 infection can also be estimated then required drug levels for protection may be identified. Here, we  
350 have demonstrated estimation of HIV infection time using non-linear viral dynamics. Specifically, we  
351 assess the viral load trajectory, an established endpoint for HIV trials. We developed a two-step  
352 procedure: 1) using population non-linear mixed effects (pNLME) modeling to estimate parameters for  
353 individuals who were infected with HIV and then 2) using these same parameters to repeatedly simulate  
354 stochastic infections. In step one, we estimate the time between first detected positive viral load and  
355 the viral load reaching some level theoretically considered deterministic. In step two, we quantify the  
356 time of stochastic viral growth until the deterministic threshold. Combining steps 1 and 2 completes the  
357 estimate of the time between exposure/acquisition and viral load detectability, sometimes called the  
358 “eclipse time”.

359 We applied our method to data from the RV217 observational cohort, an acute HIV infection study with  
360 highly granular measures of viral load early during infection. For the first step we performed extensive  
361 model selection and found a mechanistic model that recapitulates viral loads in the RV217 trial from first  
362 positive until viral set point. A trained pNLME model using data from multiple individuals that includes  
363 observations during several stages of viral infection (e.g. expansion, peak and set point) allows  
364 confidence in parameter estimation from individuals who may not have data from all stages.

365 We compared our technique to other techniques that were applied to the same data set. We find that  
366 the individual level estimates are not concordant. On the population level, average values of our  
367 deterministic model agree with average values from other approaches. The additional stochastic phase  
368 in our model drives our estimates slightly farther from time of first positive, extending the range of the  
369 eclipse time. Concordance of our model is strongest (CCC=0.2-0.4) with other approaches that use viral  
370 load dynamics. Sequence-based approaches are the least concordant, and self-report diaries are  
371 somewhat middling. We note that the study group remains was too small to evaluate certain variables  
372 that differed across individuals, such as viral subtype, sex, age or ethnicity. In cases where viral dynamics  
373 and sequencing data exist, it may be optimal to try all approaches and developing uncertainty intervals  
374 extending across all methods. A future solution would be to include evolution into our mechanistic  
375 model and fit to both types of data.

376 Compared to other methods, our approach has several advantages. It allows estimation of infection  
377 time in individuals without well resolved viral upslope or even viral peaks. Specifically, without  
378 incorporating the population data, it is not possible to estimate infection time when viral upslope is  
379 missed. It is also relatively insensitive to founder multiplicity, a challenge for phylogenetic methods that  
380 sometimes results in unrealistic infection time estimates after the first positive viral load. The RV217  
381 study is unlikely to be repeated. Thus, our model can be considered a trained model for future trials.  
382 Moreover, our synthetic data sampling study illustrates what such trials might look like.

383 While the true time of infection cannot be known other than in challenge experiments. We verified that  
384 our RV217-trained pNLME model works on simulated data (from the same mechanistic model). Even  
385 given a sparse sampling scheme (0,2,4,8,12,24 weeks after first positive) as in the antibody mediated  
386 prevention (AMP) studies, the approach generally works well for individuals for whom infection  
387 estimates are less than 20 days before first positive. This means that protocols sampling with ~2-3 week  
388 intervals typically have <20% error, or at worst 7 days off. However, our uncertainty grows for infections  
389 occurring further before first positive. This reflects the challenge of estimating data such as  
390 sparse/sparse in **Fig 7A**. Collecting only setpoint, or partial downslope means the estimate relies heavily

391 on the population model, and with the heterogeneity of individuals, can be relatively error prone.  
392 Tighter sampling after first positive does not drastically improve accuracy. Indeed, for trial design,  
393 accuracy is better-enhanced by tighter sampling prior to diagnosis.

394 The largest challenge for the approach occurs when individuals are infected at a study visit (these occur  
395 every 4 weeks) but are not diagnosed because their viral loads are below detection, meaning that the  
396 first positive viral load will be not detected until >4 weeks after infection. This challenge is not unique to  
397 our approach, and we stress diagnostic-focused assays such as APTIMA to make diagnoses as soon as  
398 possible in situations where infection timing is crucial. Thus, we have shown that 1) it is possible to  
399 leverage and impute viral loads based on the finding that early APTIMA measurements are correlated to  
400 qPCR measurements, and 2) that borrowing strength using population modeling may be the best option  
401 to overcome sparse sampling.

402 There are several limitations to our study. It remains unknown, and will be extremely hard to test,  
403 whether early HIV dynamics can be described by the same mechanistic model as deterministic viral  
404 dynamics. However, in CMV transmission the probability of infection has been related to post-infection  
405 viral kinetics, suggesting stochastic behaviors may be linked to subsequent deterministic kinetics<sup>22</sup>. We  
406 speculate an early lag period in HIV infection that could be described by localized exposure and viral  
407 escape from anatomical barriers before initiating systemic infection. The duration of this period is  
408 unknown and we effectively assumed that it is negligible compared to the stochastic and deterministic  
409 phases, and compared to our window of estimation (i.e., less than a day). To account for this lag period,  
410 a further non-mechanistic window might also be added to reconcile wider estimates of eclipse time  
411 found in some studies<sup>6</sup>. Of note, it is not clear if any timing method can directly account for this period;  
412 for example, the founder sequence may describe the virion that escapes the early barriers in our  
413 schematic.

414 Our choice of the initial simulation conditions  $I(0)$  inversely correlates with time of infection. That is, if  
415 we assume viral infection begins at a lower level, our estimates are further back in time. However, in  
416 what we consider to be a plausible range of initial conditions (ranging from 1 to 1000 infected cells  
417 initiating infection), the estimation varies by ~5 days. Interestingly, Rolland et al. found that if using a  
418 log-linear upslope modeling approach, a viral load of 1 copy/mL gave the best estimates in non-human  
419 primate infection where the date of infection was known perfectly<sup>3</sup>. One might therefore choose this  
420 value for the deterministic threshold, but the translation of this estimate is limited by the NHP  
421 experimental model, challenge virus species, and differences from viral load exposures in human  
422 transmission.

423 This approach should be generally applicable to other viruses. For example in Hepatitis C mechanistic  
424 models have been developed and some prior parameter estimates have been recorded<sup>31</sup>. Recently a  
425 similar method was applied to estimate the time of SARS-COV-2 infection<sup>32</sup>.

426 In future work, we plan to explore modifications due to preventative interventions, such that timing  
427 estimation, and therefore drug efficacy can be better estimated in upcoming clinical trials using broadly  
428 neutralizing antibodies.

429 **Acknowledgements:** DBR thanks Paul Edlefsen for motivating this work, as well as Raabya Rossenkhan,  
430 Philip Labuschagne, Dobromir Dimitrov, James Moore, Holly Janes, and Yunda Huang for help and  
431 valuable conversations. DBR is supported by the Washington Research Foundation and an NIH K25.

432

## 433 Methods

434 **Most parsimonious mathematical model.** The set of ordinary differential equations for the model that  
435 is selected for this approach can be written as

$$436 \quad \partial_t S = \alpha_S - \delta_S S - \beta SV$$

$$437 \quad \partial_t I = \beta SV - \kappa I^{h+1} \quad \text{Eq. 1}$$

$$438 \quad \partial_t V = \pi I - \gamma V - \beta SV.$$

439 The selected model contains 8 free parameters  $\theta = (\alpha_S, \delta_S, \beta, k, h, \pi, \gamma, t_{det})$ . The model we ultimately  
440 select is a slightly modified basic viral dynamics model that incorporates a nonlinear death term. The  
441 model tracks the concentration [cells mL<sup>-1</sup>] of HIV-susceptible cells  $S$ , infected cells  $I$ , and plasma viral  
442 load  $V$  [viral RNA copies mL<sup>-1</sup>]. The deterministic system is expressed (using the partial  $\partial_t$  to denote  
443 derivative in time) with  $\alpha_S$  [cells  $\mu\text{L}^{-1}$  day<sup>-1</sup>] the constant growth rate of susceptible cells,  $\delta_S$  [day<sup>-1</sup>] the  
444 death rate of susceptible cells, and  $\beta$  [ $\mu\text{L}$  virus<sup>-1</sup> day<sup>-1</sup>] a mass-action viral infectivity. The viral production  
445 rate is  $\pi$  [virions cells<sup>-1</sup> day<sup>-1</sup>], and  $\gamma$  [day<sup>-1</sup>] is the clearance rate of virus. The death and killing of infected  
446 cells is governed by the rate of  $\kappa$  [cells<sup>-h</sup> day<sup>-1</sup>], with the exponential factor  $h$  adjusting the nonlinear  
447 density dependent death rate. This approach coarsely approximates adaptive immunity such that higher  
448 numbers of infected cells engenders faster killing.

449 **Population nonlinear mixed effects (pNLME) approach.** We modeled the plasma viral load using a  
450 nonlinear mixed-effects approach (pNLME). In this approach an observed plasma viral load for individual  
451  $i$  at time  $j$  is modeled as  $\log_{10} y_{ij} = f_V(t_{ij}, \theta_i) + \epsilon_V$ . Here,  $f_V$  is the solution of the nonlinear  
452 mechanistic model for the variable describing the virus ( $V$ ) given the individual parameter vector  $\theta_i$  and  
453  $\epsilon_V \sim \mathcal{N}(0, \sigma_v^2)$  is the measurement error for the logged viral load. We assumed that the individual-  
454 specific parameter  $\theta_i$  is drawn from a probability distribution with median or fixed effects  $\theta^{pop}$  and  
455 random effects  $\eta_i \sim \mathcal{N}(0, \Omega)$ , being  $\Omega$  the variance-covariance matrix. Except otherwise specified we  
456 modeled parameters  $\beta^j$  and  $\pi^j$  as  $\theta_i = 10^{\theta^{pop} + \eta_i}$  and remaining parameters as  $\theta_j = \theta^{pop} e^{\eta_j}$ .

457  
458 **Model fitting.** We explored four different mechanistic models with different statistical complexities, for a  
459 total of 30 models (See **Supplementary Table 1** for details). For each model we obtained the Maximum  
460 Likelihood Estimation (MLE) of the measurement error standard deviation  $\sigma_v$ , the fixed effects vector  
461  $\theta^{pop}$  and the elements of matrix  $\Omega$  using the Stochastic Approximation of the Expectation Maximization  
462 (SAEM) algorithm embedded in the Monolix software ([www.lixoft.eu](http://www.lixoft.eu)). We run the SAEM algorithm 15  
463 times (assessments) for each model using randomly selected initial guesses for the parameters to  
464 estimate. For all model fits we assumed  $t_{ij} = 0$  as the time of first positive viral load. However, we  
465 defined the initial value as the time  $-t_{det}$  when  $V(-t_{det}) = 0.01$  copies/mL. We fixed other initial values  
466 as  $S(t_{det}) = \frac{\alpha_S}{\delta_S}$  cells/ $\mu\text{L}$  and  $I(t_{det}) = \frac{\gamma V(-t_{det})}{\pi}$  cells/ $\mu\text{L}$ . Per Ref <sup>31</sup>, we fixed parameter  $\gamma = 23$  day<sup>-1</sup>. We  
467 estimated the remaining parameters of the mechanistic model including  $t_{det}$ . Individual parameters  
468 were selected using the mode of the conditional distribution  $p(\theta_i | y_{ij}; \theta_{MLE}^{pop}, \Omega_{MLE})$  constructed by the  
469 MCMC algorithm in the Monolix software. The conditional distribution of  $t_{det}$  for each individual is used  
470 to compute the time of infection  $t_0$ .

471  
472 **Model selection.** To determine the most parsimonious model we calculated the log-likelihood ( $\log L$ ) for  
473 all 15 assessments for each one of the 30 models. We then computed the Akaike Information Criteria for  
474 the model with highest likelihood among the 15 assessments ( $AIC = -2 \log \mathcal{L}_{max} + 2m$ , where  $m$  is  
475 the number of parameters estimated). We assumed a model has similar support from the data if the

476 difference between the AIC for its best assessment and the best one for the model with lowest AIC is  
477 less than two<sup>33</sup>.

478  
479 **Stochastic simulation scheme.** We adapted the ordinary differential equation system Eq 1 to simulate  
480 stochastically<sup>12</sup>. Our implementation in Python, which employs the  $\tau$ -leap approach<sup>24</sup>, is publicly  
481 available. A time interval  $\Delta t = 0.0001$  days is chosen for step size, in which a Poisson number of each  
482 transition type occurs. Initial conditions are changed to discrete values by multiplying by a volume. We  
483 choose this volume to be  $10^8$   $\mu\text{L}$  based on the observation that there is approximately 1-10 L of blood in  
484 an adult human and that there are approximately 10- 100 times more T cells in lymph tissue than blood.  
485 A single infected cell is assumed (other than in sensitivity analyses in **Supplementary figure 4**).

486 **APTIMA analysis.** APTIMA was the primary diagnostic assay in the RV217 study, of which 43  
487 participants in our analysis were diagnosed by a positive APTIMA quantitative measurement. The  
488 remaining participants were diagnosed via a qualitative APTIMA response or directly with viral load.  
489 Among the 43 participants, only 6 had concurrent measurements of viral load for analysis. Comparing  
490 concurrent measurements APTIMA and viral load we found 1) substantial Pearson correlation between  
491 APTIMA and viral load at first positive viral load; and 2) diminished correlation later in the study at  
492 higher values of both measurements (**Figure 2A**).

493 Given the high correlation between the two measurements, we sought to investigate whether we could  
494 use APTIMA measurements at diagnosis to predict the unmeasured viral load for our model. This was  
495 accomplished using linear regression models predicting log first positive viral load with concurrent  
496 APTIMA as the predictor, evaluating both untransformed and log-transformed inputs (**Supplementary**  
497 **figure 6A**). One participant had an APTIMA measurement of 3, an outlier more than 2-fold lower than  
498 the next lowest value, and were removed from the model. To determine the appropriate upper range  
499 for APTIMA input for prediction, linear regression models were fit applying different upper bounds.  
500 Model performance was evaluated using residual mean square error (RMSE) predicting log viral load  
501 (**Supplementary figure 6B**). We found the best model used raw APTIMA measurements as the input  
502 with an upper bound of 34 (**Supplementary figure 6B&C**). This model was applied to the data to impute  
503 first positives for participants where APTIMA was measured for diagnosis without viral load  
504 (**Supplementary figure 6D**).

## 505 References

- 506  
507 1 Giorgi EE, Funkhouser B, Athreya G, Perelson AS, Korber BT, Bhattacharya T. Estimating time  
508 since infection in early homogeneous HIV-1 samples using a poisson model. *BMC Bioinformatics*  
509 2010;**11**:532. <https://doi.org/10.1186/1471-2105-11-532>.
- 510 2 Puller V, Neher R, Albert J. Estimating time of HIV-1 infection from next-generation sequence  
511 diversity. *PLoS Comput Biol* 2017;**13**:1–20. <https://doi.org/10.1371/journal.pcbi.1005775>.
- 512 3 Rolland M, Tovanabutra S, Dearlove B, Li Y, Owen CL, Lewitus E, *et al.* Molecular dating and viral  
513 load growth rates suggested that the eclipse phase lasted about a week in HIV-1 infected adults  
514 in East Africa and Thailand. *PLOS Pathog* 2020;**16**:e1008179.  
515 <https://doi.org/10.1371/journal.ppat.1008179>.
- 516 4 Fiebig EW, Busch MP, Wright DJ, Kleinman SH, Rawal BD, Garrett PE, *et al.* Dynamics of HIV  
517 viremia and antibody seroconversion in plasma donors: Implications for diagnosis and staging of  
518 primary HIV infection. *Aids* 2003;**17**:1871–9.  
519 <https://doi.org/10.1097/01.aids.0000076308.76477.b8>.

- 520 5 Delaney KP, Hanson DL, Masciotra S, Ethridge SF, Wesolowski L, Owen SM. Time Until Emergence  
521 of HIV Test Reactivity Following Infection With HIV-1 : Implications for Interpreting Test Results  
522 and Retesting After Exposure 2017;**64**:53–9. <https://doi.org/10.1093/cid/ciw666>.
- 523 6 Pilcher CD, Porco TC, Facente SN, Grebe E, Delaney KP, Masciotra S, *et al.* A generalizable method  
524 for estimating duration of HIV infections using clinical testing history and HIV test results. *AIDS*  
525 2019;**33**:1231–40. <https://doi.org/10.1097/QAD.0000000000002190>.
- 526 7 Rossenkhan R, Rolland M, Labuschagne JPL, Ferreira R, Magaret CA, Carpp LN, *et al.* Combining  
527 Viral Genetics and Statistical Modeling to Improve HIV-1 Time-of-infection Estimation towards  
528 Enhanced Vaccine E ffi cacy Assessment 2019.
- 529 8 Robb ML, Eller LA, Kibuuka H, Rono K, Maganga L, Nitayaphan S, *et al.* Prospective Study of Acute  
530 HIV-1 Infection in Adults in East Africa and Thailand. *N Engl J Med* 2016;NEJMoa1508952.  
531 <https://doi.org/10.1056/NEJMoa1508952>.
- 532 9 Cohen MS, Shaw GM, McMichael AJ, Haynes BF. Acute HIV-1 Infection. *N Engl J Med*  
533 2011;**364**:1943–54. <https://doi.org/10.1056/NEJMra1011874>.
- 534 10 Konrad BP, Taylor D, Conway JM, Ogilvie GS, Coombs D. On the duration of the period between  
535 exposure to HIV and detectable infection. *Epidemics* 2017;**20**:73–83.  
536 <https://doi.org/10.1016/j.epidem.2017.03.002>.
- 537 11 Conway JM, Konrad BP, Coombs D. Stochastic Analysis of Pre- and Postexposure Prophylaxis  
538 against HIV Infection. *SIAM J Appl Math* 2013;**73**:904–28. <https://doi.org/10.1137/120876800>.
- 539 12 Reeves DB, Huang Y, Duke ER, Mayer BT, Fabian Cardozo-Ojeda E, Boshier FA, *et al.* Mathematical  
540 modeling to reveal breakthrough mechanisms in the HIV Antibody Mediated Prevention (AMP)  
541 trials. *PLoS Comput Biol* 2020;**16**:1–27. <https://doi.org/10.1371/journal.pcbi.1007626>.
- 542 13 Liu J, Ghneim K, Sok D, Bosche WJ, Li Y, Chipriano E, *et al.* Antibody-mediated protection against  
543 SHIV challenge includes systemic clearance of distal virus. *Science (80- )* 2016;**353**:1045–9.  
544 <https://doi.org/10.1126/science.aag0491>.
- 545 14 Hessel AJ, Pognard P, Hunter M, Hangartner L, Tehrani DM, Bleeker WK, *et al.* Effective, low-  
546 titer antibody protection against low-dose repeated mucosal SHIV challenge in macaques. *Nat*  
547 *Med* 2009;**15**:951–4. <https://doi.org/10.1038/nm.1974>.
- 548 15 Perelson AS, Ribeiro RM. Modeling the within-host dynamics of HIV infection. *BMC Biol*  
549 2013;**11**:96. <https://doi.org/10.1186/1741-7007-11-96>.
- 550 16 Reeves DB, Peterson CW, Kiem H, Schiffer JT. Autologous Stem Cell Transplantation Disrupts  
551 Adaptive Immune Responses during Rebound Simian/Human Immunodeficiency Virus Viremia. *J*  
552 *Viro* 2017;**91**:e00095--17. <https://doi.org/10.1128/JVI.00095-17>.
- 553 17 Borducchi EN, Liu J, Nkolola JP, Cadena AM, Yu W-H, Fischinger S, *et al.* Antibody and TLR7  
554 agonist delay viral rebound in SHIV-infected monkeys. *Nature* 2018.  
555 <https://doi.org/10.1038/s41586-018-0600-6>.
- 556 18 Holte SE, Melvin AJ, Mullins JI, Tobin NH, Frenkel LM. Density-dependent decay in HIV-1  
557 dynamics. *JAIDS J Acquir Immune Defic Syndr* 2006;**41**:266–76.  
558 <https://doi.org/10.1097/01.qai.0000199233.69457.e4>.
- 559 19 Davenport MP, Zhang L, Shiver JW, Casmiro DR, Ribeiro RM, Perelson AS. Influence of peak viral

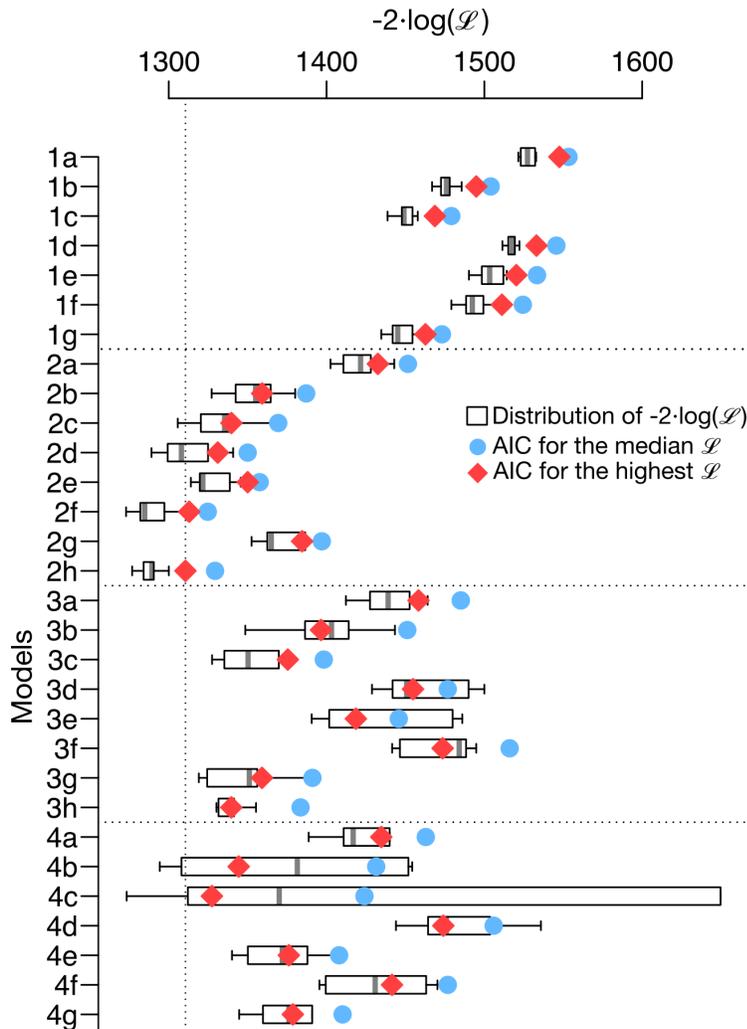
- 560 load on the extent of CD4+ T-cell depletion in simian HIV infection. *J Acquir Immune Defic Syndr*  
561 2006;**41**:259–65. [https://doi.org/10.1097/01.qai.0000199232.31340.d3.n00126334-200603000-](https://doi.org/10.1097/01.qai.0000199232.31340.d3.n00126334-200603000-00001)  
562 00001 [pii].
- 563 20 Robb ML, Eller LA, Kibuuka H, Rono K, Maganga L, Nitayaphan S, *et al.* Prospective Study of Acute  
564 HIV-1 Infection in Adults in East Africa and Thailand. *N Engl J Med* 2016;NEJMoa1508952.  
565 <https://doi.org/10.1056/NEJMoa1508952>.
- 566 21 Mayer BT, Krantz EM, Swan D, Ferrenberg J, Simmons K, Selke S, *et al.* Transient Oral Human  
567 Cytomegalovirus Infections Indicate Inefficient Viral Spread from Very Few Initially Infected Cells.  
568 *J Virol* 2017;**91**:2701–13. <https://doi.org/10.1128/JVI.00380-17>.
- 569 22 Mayer BT, Matrajt L, Casper C, Krantz EM, Corey L, Wald A, *et al.* Dynamics of persistent oral  
570 cytomegalovirus shedding during primary infection in ugandan infants. *J Infect Dis*  
571 2016;**214**:1735–43. <https://doi.org/10.1093/infdis/jiw442>.
- 572 23 Ganusov V V, De Boer RJ. Do most lymphocytes in humans really reside in the gut? *Trends*  
573 *Immunol* 2007;**28**:514–8. <https://doi.org/10.1016/j.it.2007.08.009>.
- 574 24 Gillespie DT. Approximate accelerated stochastic simulation of chemically reacting systems. *J*  
575 *Chem Phys* 2001;**115**:1716–33. <https://doi.org/10.1063/1.1378322>.
- 576 25 Hughes JP, Baeten JM, Lingappa JR, Magaret AS, Wald A, de Bruyn G, *et al.* Determinants of Per-  
577 Coital-Act HIV-1 Infectivity Among African HIV-1–Serodiscordant Couples. *J Infect Dis*  
578 2012;**205**:358–65. <https://doi.org/10.1093/infdis/jir747>.
- 579 26 Drummond AJ, Suchard MA, Xie D, Rambaut A. Bayesian phylogenetics with BEAUti and the  
580 BEAST 1.7. *Mol Biol Evol* 2012;**29**:1969–73. <https://doi.org/10.1093/molbev/mss075>.
- 581 27 Giorgi EE, Funkhouser B, Athreya G, Perelson AS, Korber BT, Bhattacharya T. Estimating time  
582 since infection in early homogeneous HIV-1 samples using a poisson model. *BMC Bioinformatics*  
583 2010;**11**:532. <https://doi.org/10.1186/1471-2105-11-532>.
- 584 28 Janes H, Herbeck JT, Tovanabutra S, Thomas R, Frahm N, Duerr A, *et al.* HIV-1 infections with  
585 multiple founders are associated with higher viral loads than infections with single founders. *Nat*  
586 *Med* 2015;**21**:1139–41. <https://doi.org/10.1038/nm.3932>.
- 587 29 Gilbert PB, Juraska M, DeCamp AC, Karuna S, Edupuganti S, Mgodhi N, *et al.* Basis and Statistical  
588 Design of the Passive HIV-1 Antibody Mediated Prevention (AMP) Test-of-Concept Efficacy Trials.  
589 *Stat Commun Infect Dis* 2017;**9**:. <https://doi.org/10.1515/scid-2016-0001>.
- 590 30 Huang Y, Karuna S, Carpp LN, Reeves D, Pegu A, Seaton K, *et al.* Modeling cumulative overall  
591 prevention efficacy for the VRC01 phase 2b efficacy trials. *Hum Vaccines Immunother* 2018;**0**:1–  
592 12. <https://doi.org/10.1080/21645515.2018.1462640>.
- 593 31 Ramratnam B, Bonhoeffer S, Binley J, Hurley A, Zhang L, Mittler JE, *et al.* Rapid production and  
594 clearance of HIV-1 and hepatitis C virus assessed by large volume plasma apheresis. *Lancet*  
595 1999;**354**:1782–5. [https://doi.org/10.1016/S0140-6736\(99\)02035-8](https://doi.org/10.1016/S0140-6736(99)02035-8).
- 596 32 Ejima K, Kim KS, Ito Y, Iwanami S, Ohashi H, Koizumi Y, *et al.* Inferring Timing of Infection Using  
597 Within-host SARS-CoV-2 Infection Dynamics Model: Are ‘Imported Cases’ Truly Imported?  
598 *MedRxiv* 2020;**4297**:2020.03.30.20040519. <https://doi.org/10.1101/2020.03.30.20040519>.
- 599 33 Beier P, Burnham KP, Anderson DR. *Model Selection and Inference: A Practical Information-*

600 *Theoretic Approach*. vol. 65. 2001.

601

602 **Supplementary figures and tables**

603

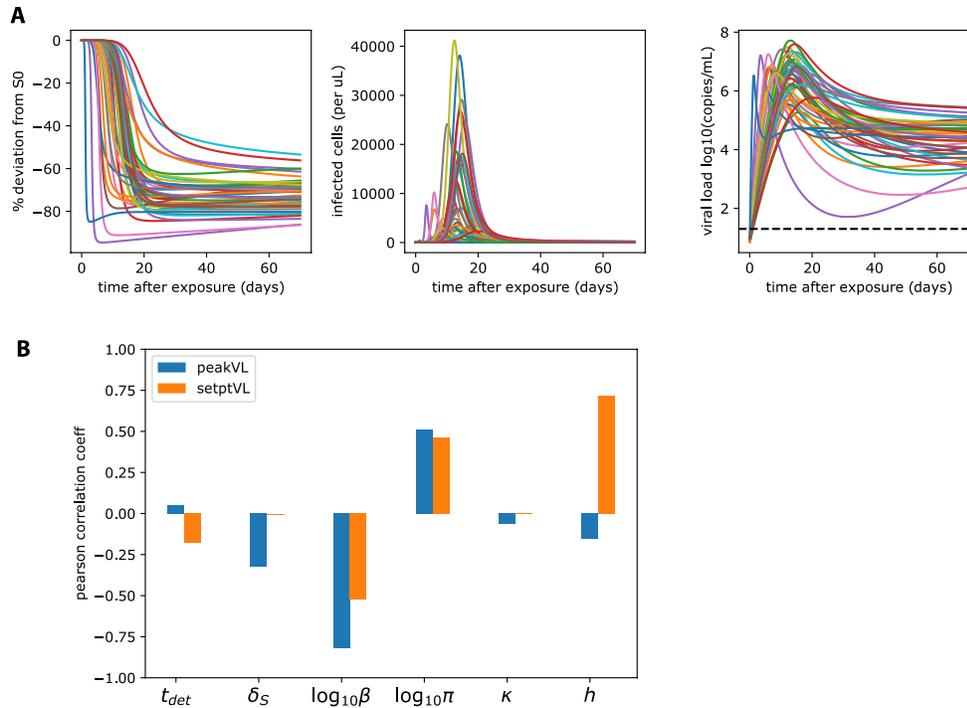


604

605 **Supplementary figure 1. Model selection details.** All 30 models tested compared by log likelihood (-2LL)  
606 and Akaike Information Criterion (AIC). Boxplots give -2LL of 15 assessments for each model, where each  
607 assessment begins at a different parameter set and proceeds with the stochastic SAEM algorithm. Red  
608 diamonds give AIC of the median -2LL across assessments. Blue circles give the AIC of the mean  
609 assessment. The best model is 2h. Several different correlation models were attempted for each model,  
610 ultimately leading to the most identifiable model combinations.

611

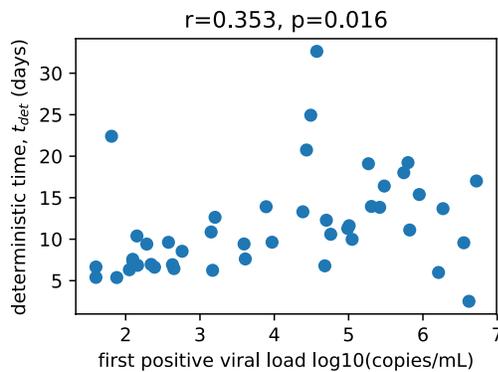
It is made available under a [CC-BY-ND 4.0 International license](https://creativecommons.org/licenses/by-nd/4.0/).



612

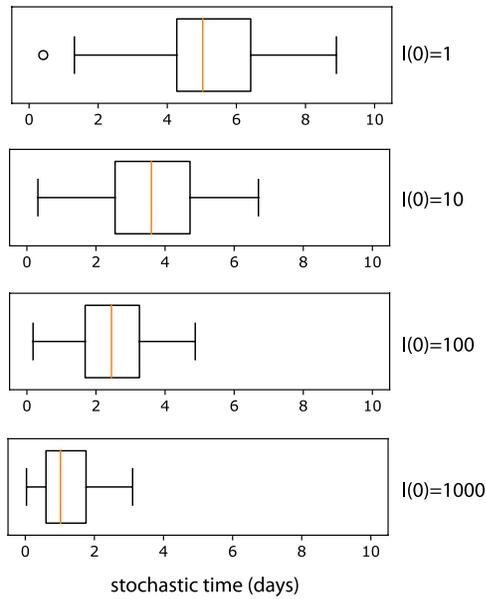
613 **Supplementary figure 2. Model sensitivity analysis.** All best fit parameter sets simulated together. A) %  
 614 deviation from the initial number of susceptible cells can go up to -100% percent, indicating massive  
 615 destruction of cells in acute HIV infection. The total number of infected cells at that point can rise to  
 616 ~1000 cells per  $\mu\text{L}$ . Viral loads can have peaks ranging from  $10^6$ - $10^8$  copies/mL, with setpoints varying  
 617 substantially between  $10^2$ - $10^6$  copies/mL. B) Pearson correlation between parameters and viral kinetic  
 618 phenotypes.

619



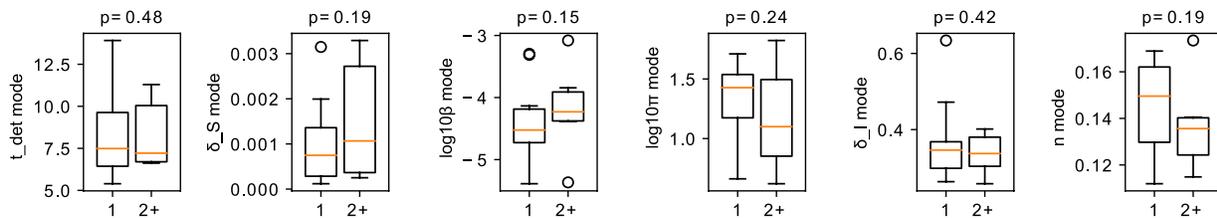
620

621 **Supplementary figure 3.** Correlation between observed first positive viral load and the inferred  
 622 deterministic time.



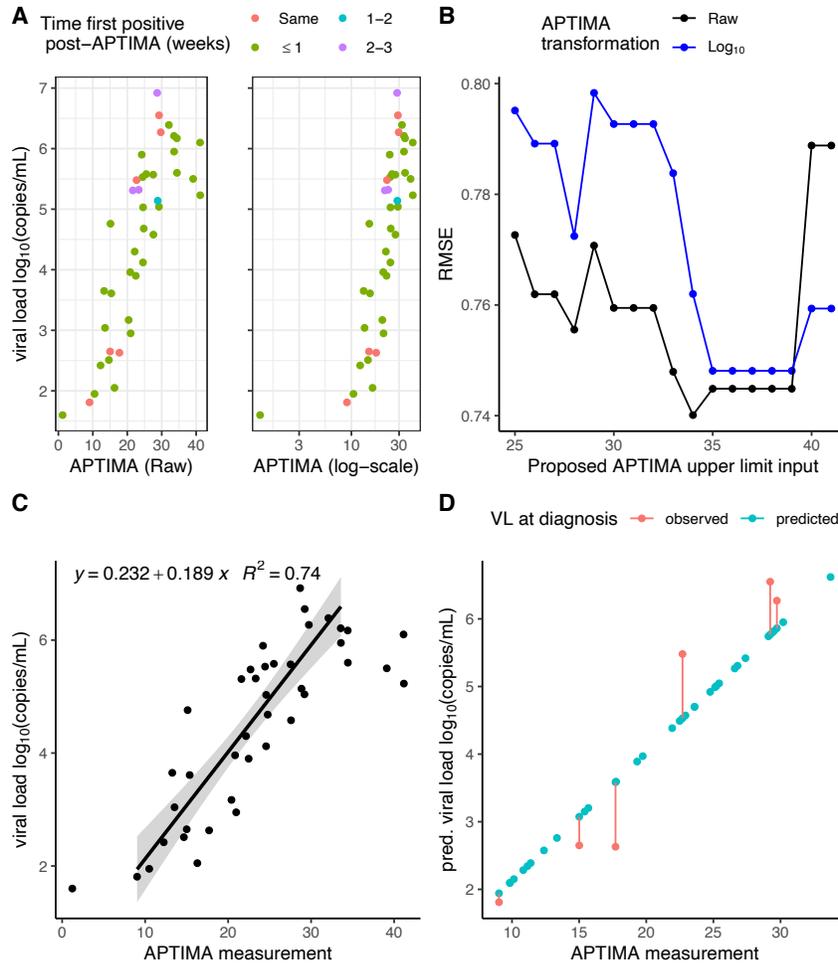
623

624 **Supplementary figure 4. Relationship between initial number of infected cells and stochastic time.** As  
625 the initial number of infected cells is increased, the time to reach the deterministic threshold decreases.  
626 Within this plausible range of 1-1000 initially infected cells, the estimates of the stochastic interval  
627 decrease from median 5 to median 1 day. This affects estimations, bounded by this 4 day window. For  
628 each value of  $I(0)$ , 1000 simulations were performed.



629

630 **Supplementary figure 5. Relationship between viral dynamic parameters and multiple founder status.**  
631 Grouped as either 1 founder ( $n = 24$ ) or 2+ founders ( $n = 9$ ) No estimated parameter is significantly  
632 (Mann-Whitney U-test  $p$ -value above each panel) different across founders, meaning the method is not  
633 significantly affected by the presence of multiple founders.



634

635 **Supplementary figure 6. Predicting viral load (VL) with APTIMA measurement at diagnosis.** A) First  
 636 positive viral load measurements vs. concurrent raw or log-transformed APTIMA measurements. Color  
 637 denotes time (weeks) of first positive relative to the diagnostic APTIMA measurement. 'Same' denotes  
 638 participants who were diagnosed by first positive viral load and positive APTIMA on the same visit. B)  
 639 Residual mean squared error (RMSE) predicting  $\log_{10}$  first positive viral load with concurrent APTIMA  
 640 (raw or log-transformed) for varying input data for different APTIMA upper bounds. C) Selected best  
 641 regression model from B) denoted by the line with shading for standard error for predicting viral load  
 642 where APTIMA input range limited to 9-34. Raw data denoted by points. D) Predicted first positive viral  
 643 load using the model depicted in C) and participants' APTIMA measurements at diagnosis. Red dots  
 644 denote the 6 participants where first positive and diagnostic APTIMA were measured together, and red  
 645 line depicts prediction error. Predicted viral load was only used when for participants without viral load  
 646 measurements at diagnosis.

647 **Supplementary Table 1.** Description of all the models with their assumptions that were fit to the data.  
 648 Estimated standard deviation of random effects ( $\sigma$ ) and correlations for matrix  $\Omega$  for each model are  
 649 specified, if not included in the table they were assumed to be zero. Other fixed parameters not specified  
 650 in the text are included here. Distribution for the Log likelihood estimations for each model, AIC for the  
 651 median and highest likelihood are presented in **Supplementary figure 1**.

652

Model	Model Name	Fixed Parameters	Estimated Random Effects	Estimated Correlations
<b>1. Basic:</b> $dS/dt = \alpha_S - \delta_S S - \beta SV$ ; $dI/dt = \beta SV - dI$ ; $dV/dt = \pi I - \gamma V - \beta SV$	<b>1a</b>	No other fixed parameters	$\sigma_{tdet}, \sigma_{\alpha_S}, \sigma_{\delta_S}, \sigma_{\log10\beta}, \sigma_{\log10\pi}, \sigma_{\delta I}$	No correlations
	<b>1b</b>	No other fixed parameters	$\sigma_{tdet}, \sigma_{\alpha_S}, \sigma_{\delta_S}, \sigma_{\log10\beta}, \sigma_{\log10\pi}, \sigma_{\delta I}$	$\text{corr}(\log10\beta, \alpha_S)$
	<b>1c</b>	No other fixed parameters	$\sigma_{tdet}, \sigma_{\alpha_S}, \sigma_{\delta_S}, \sigma_{\log10\beta}, \sigma_{\log10\pi}, \sigma_{\delta I}$	$\text{corr}(tdet, \delta_S), \text{corr}(\log10\beta, \alpha_S)$
	<b>1d</b>	No other fixed parameters	$\sigma_{tdet}, \sigma_{\alpha_S}, \sigma_{\delta_S}, \sigma_{\log10\beta}, \sigma_{\log10\pi}, \sigma_{\delta I}$	$\text{corr}(\log10\pi, \log10\beta)$
	<b>1e</b>	No other fixed parameters	$\sigma_{tdet}, \sigma_{\alpha_S}, \sigma_{\delta_S}, \sigma_{\log10\beta}, \sigma_{\log10\pi}, \sigma_{\delta I}$	$\text{corr}(\log10\beta, \delta_S), \text{corr}(\log10\pi, dI)$
	<b>1f</b>	No other fixed parameters	$\sigma_{tdet}, \sigma_{\alpha_S}, \sigma_{\delta_S}, \sigma_{\log10\beta}, \sigma_{\log10\pi}, \sigma_{\delta I}$	$\text{corr}(\log10\beta, \delta_S), \text{corr}(\log10\pi, \delta_S), \text{corr}(\log10\pi, \log10\beta)$
	<b>1g</b>	No other fixed parameters	$\sigma_{tdet}, \sigma_{\alpha_S}, \sigma_{\delta_S}, \sigma_{\log10\beta}, \sigma_{\delta I}$	$\text{corr}(tdet, \delta_S), \text{corr}(\log10\beta, \alpha_S)$
<b>2. Holte:</b> $dS/dt = \alpha_S - \delta_S S - \beta SV$ ; $dI/dt = \beta SV - \kappa I^{h+1}$ ; $dV/dt = \pi I - \gamma V - \beta SV$	<b>2a</b>	No other fixed parameters	$\sigma_{tdet}, \sigma_{\alpha_S}, \sigma_{\delta_S}, \sigma_{\log10\beta}, \sigma_{\log10\pi}, \sigma_{\kappa}, \sigma_h$	No correlations
	<b>2b</b>	No other fixed parameters	$\sigma_{tdet}, \sigma_{\alpha_S}, \sigma_{\delta_S}, \sigma_{\log10\beta}, \sigma_{\log10\pi}, \sigma_{\kappa}, \sigma_h$	$\text{corr}(\log10\pi, \log10\beta)$
	<b>2c</b>	No other fixed parameters	$\sigma_{tdet}, \sigma_{\alpha_S}, \sigma_{\delta_S}, \sigma_{\log10\beta}, \sigma_{\log10\pi}, \sigma_{\kappa}, \sigma_h$	$\text{corr}(\log10\beta, \delta_S), \text{corr}(\log10\pi, dI)$
	<b>2d</b>	No other fixed parameters	$\sigma_{tdet}, \sigma_{\alpha_S}, \sigma_{\delta_S}, \sigma_{\log10\beta}, \sigma_{\log10\pi}, \sigma_{\kappa}, \sigma_h$	$\text{corr}(\delta_S, dI), \text{corr}(\log10\beta, dI), \text{corr}(\log10\pi, dI), \text{corr}(\log10\beta, \delta_S), \text{corr}(\log10\pi, \delta_S), \text{corr}(\log10\pi, \log10\beta)$
	<b>2e</b>	No other fixed parameters	$\sigma_{tdet}, \sigma_{\alpha_S}, \sigma_{\delta_S}, \sigma_{\log10\beta}, \sigma_{\log10\pi}, \sigma_{\kappa}, \sigma_h$	$\text{corr}(\log10\beta, \delta_S), \text{corr}(\log10\pi, \delta_S), \text{corr}(\log10\pi, \log10\beta)$
	<b>2f</b>	No other fixed parameters	$\sigma_{tdet}, \sigma_{\delta_S}, \sigma_{\log10\beta}, \sigma_{\log10\pi}, \sigma_{\kappa}, \sigma_h$	$\text{corr}(\delta_S, dI), \text{corr}(\log10\beta, dI), \text{corr}(\log10\pi, dI), \text{corr}(\log10\beta, \delta_S), \text{corr}(\log10\pi, \delta_S), \text{corr}(\log10\pi, \log10\beta)$
	<b>2g</b>	No other fixed parameters	$\sigma_{tdet}, \sigma_{\delta_S}, \sigma_{\log10\beta}, \sigma_{\log10\pi}, \sigma_h$	$\text{corr}(\log10\beta_{\delta_S}, \text{corr}(\log10\pi_{\delta_S}, \text{corr}(\log10\pi, \log10\beta)$
	<b>2h</b>	No other fixed parameters	$\sigma_{tdet}, \sigma_{\delta_S}, \sigma_{\log10\beta}, \sigma_{\log10\pi}, \sigma_{\kappa}, \sigma_h$	$\text{corr}(\delta_S, dI), \text{corr}(\log10\beta, dI), \text{corr}(\log10\pi, dI), \text{corr}(\log10\beta, \delta_S), \text{corr}(\log10\pi, \delta_S), \text{corr}(\log10\pi, \log10\beta)$

<p><b>3. Hill:</b>  <math>dS/dt = \alpha_S - \delta_S S - \beta SV / (1 + \phi E)</math>;  <math>dI/dt = \beta SV / (1 + \phi E) - \delta_I I - \kappa EI</math>;  <math>dV/dt = \pi I - \gamma V - \beta SV</math>;  <math>dP/dt = \alpha_E + \omega(1-f)PI / (1 + I/N) - \delta_P P</math>;  <math>dE/dt = \omega fPI / (1 + I/N) - \delta_E E</math></p>	<b>3a</b>	$\delta E_{pop}=1$	$\sigma_{tdet}, \sigma_{\alpha_S}, \sigma_{\delta_S}, \sigma_{\log 10 \beta}, \sigma_{\log 10 \pi}, \sigma_{\delta_I}, \sigma_{\log 10 \phi}, \sigma_{\alpha_E}, \sigma_{\omega}, \sigma_h, \sigma_{\delta_P}$	No correlations
	<b>3b</b>	$\delta E_{pop}=1$	$\sigma_{tdet}, \sigma_{\alpha_S}, \sigma_{\delta_S}, \sigma_{\log 10 \beta}, \sigma_{\log 10 \pi}, \sigma_{\delta_I}, \sigma_{\log 10 \phi}, \sigma_{\alpha_E}, \sigma_{\omega}, \sigma_h, \sigma_{\delta_P}$	corr(log10 $\beta$ ,tdet)
	<b>3c</b>	$\delta E_{pop}=1$	$\sigma_{tdet}, \sigma_{\alpha_S}, \sigma_{\delta_S}, \sigma_{\log 10 \beta}, \sigma_{\log 10 \pi}, \sigma_{\delta_I}, \sigma_{\log 10 \phi}, \sigma_{\alpha_E}, \sigma_{\omega}, \sigma_h, \sigma_{\delta_P}$	corr(log10 $\pi$ ,log10 $\beta$ )
	<b>3d</b>	$\delta S_{pop}=0.05, \log 10 \pi_{pop}=4.7, \delta I_{pop}=0.4, \log 10 \phi_{pop}=-4, \delta P_{pop}=0.001, \delta E_{pop}=1$	$\sigma_{tdet}, \sigma_{\alpha_S}, \sigma_{\log 10 \beta}, \sigma_{\alpha_E}, \sigma_{\omega}, \sigma_h$	No correlations
	<b>3e</b>	$\delta S_{pop}=0.05, \log 10 \pi_{pop}=4.7, \delta I_{pop}=0.4, \log 10 \phi_{pop}=-4, \delta P_{pop}=0.001, \delta E_{pop}=1$	$\sigma_{tdet}, \sigma_{\alpha_S}, \sigma_{\log 10 \beta}, \sigma_{\alpha_E}, \sigma_{\omega}, \sigma_h$	corr(tdets, $\alpha_S$ )
	<b>3f</b>	$\delta S_{pop}=0.05, \delta I_{pop}=0.4, \log 10 \phi_{pop}=-4, \delta P_{pop}=0.001, \delta E_{pop}=1$	$\sigma_{tdet}, \sigma_{\alpha_S}, \sigma_{\log 10 \beta}, \sigma_{\log 10 \pi}, \sigma_{\alpha_E}, \sigma_{\omega}, \sigma_h$	corr(log10 $\pi$ ,log10 $\beta$ )
	<b>3g</b>	log10 $\phi_{pop}=-3, \alpha E_{pop}=0.1, \delta E_{pop}=1$	$\sigma_{tdet}, \sigma_{\alpha_S}, \sigma_{\delta_S}, \sigma_{\log 10 \beta}, \sigma_{\log 10 \pi}, \sigma_{\delta_I}, \sigma_{\omega}, \sigma_h, \sigma_{\delta_P}$	corr(log10 $\pi$ ,log10 $\beta$ )
	<b>3h</b>	log10 $\phi_{pop}=-3, \alpha E_{pop}=0.1, \delta E_{pop}=1$	$\sigma_{tdet}, \sigma_{\alpha_S}, \sigma_{\delta_S}, \sigma_{\log 10 \beta}, \sigma_{\log 10 \pi}, \sigma_{\delta_I}, \sigma_{\omega}, \sigma_h, \sigma_{\delta_P}$	corr(log10 $\beta$ ,tdet), corr(log10 $\pi$ ,tdet), corr(log10 $\pi$ ,log10 $\beta$ )
<p><b>4. Reeves:</b>  <math>dS/dt = \alpha_S - \delta_S S - \beta SV</math>;  <math>dI_p/dt = \tau \beta SV - \delta_{I_p} I_p - \kappa EI_p</math>;  <math>dI_u/dt = (1-\tau) \beta SV - \delta_{I_u} I_u - \kappa EI_u</math>;  <math>dV/dt = \pi I_p - \gamma V - \beta SV</math>;  <math>dE/dt = \alpha_E + \omega EI / (E + E_{50}) - \delta_E E</math></p>	<b>4a</b>	No other fixed parameters	$\sigma_{tdet}, \sigma_{\alpha_S}, \sigma_{\delta_S}, \sigma_{\log 10 \beta}, \sigma_{\log 10 \pi}, \sigma_{\delta_I}, \sigma_{\log 10 \kappa}, \sigma_{\alpha_E}, \sigma_{\omega}, \sigma_{E50}, \sigma_{\delta_E}$	No correlations
	<b>4b</b>	No other fixed parameters	$\sigma_{tdet}, \sigma_{\alpha_S}, \sigma_{\delta_S}, \sigma_{\log 10 \beta}, \sigma_{\log 10 \pi}, \sigma_{\delta_I}, \sigma_{\log 10 \kappa}, \sigma_{\alpha_E}, \sigma_{\omega}, \sigma_{E50}, \sigma_{\delta_E}$	corr(tdets, $\alpha_S$ ), corr(log10 $\kappa$ , $\delta_S$ )
	<b>4c</b>	No other fixed parameters	$\sigma_{tdet}, \sigma_{\alpha_S}, \sigma_{\delta_S}, \sigma_{\log 10 \beta}, \sigma_{\log 10 \pi}, \sigma_{\delta_I}, \sigma_{\log 10 \kappa}, \sigma_{\alpha_E}, \sigma_{\omega}, \sigma_{E50}, \sigma_{\delta_E}$	corr( $\delta_E$ , $\alpha_S$ ), corr(tdets, $\alpha_S$ ), corr(tdets, $\delta_E$ ), corr(log10 $\kappa$ , $\delta_S$ )

	<b>4d</b>	No other fixed parameters	$\sigma_{t_{det}}, \sigma_{\delta_S}, \sigma_{\log_{10}\beta}, \sigma_{\log_{10}\kappa}, \sigma_{\omega}, \sigma_{E_{50}}, \sigma_{\delta E}$	No correlations
	<b>4e</b>	$\alpha S_{pop}=70, \log_{10}\pi_{pop}=4.7, \delta I_{pop}=0.8, \alpha E_{pop}=1e-05$	$\sigma_{t_{det}}, \sigma_{\delta_S}, \sigma_{\log_{10}\beta}, \sigma_{\log_{10}\kappa}, \sigma_{\omega}, \sigma_{E_{50}}, \sigma_{\delta E}$	corr(log10 $\beta$ , tdet), corr(log10 $\kappa$ , tdet), corr(log10 $\kappa$ , log10 $\beta$ )
	<b>4f</b>	No other fixed parameters	$\sigma_{t_{det}}, \sigma_{\alpha_S}, \sigma_{\delta_S}, \sigma_{\log_{10}\beta}, \sigma_{\log_{10}\pi}, \sigma_{\log_{10}\kappa}, \sigma_{\omega}, \sigma_{E_{50}}, \sigma_{\delta E}$	corr( $\delta E$ , $\alpha_S$ ), corr(tdet, $\alpha_S$ ), corr(tdet, $\delta E$ ), corr(log10 $\kappa$ , $\delta_S$ )
	<b>4g</b>	$\alpha S_{pop}=70, \log_{10}\pi_{pop}=4.7, \delta I_{pop}=0.8, \alpha E_{pop}=1e-05$	$\sigma_{t_{det}}, \sigma_{\delta_S}, \sigma_{\log_{10}\beta}, \sigma_{\log_{10}\kappa}, \sigma_{E_{50}}, \sigma_{\delta E}$	corr(log10 $\beta$ , tdet), corr(log10 $\kappa$ , tdet), corr(log10 $\kappa$ , log10 $\beta$ )

653

654 **Supplementary Table 2.** 6 estimated parameter estimates for each individual. Fixed parameters include  
 655  $V_0=0.01$  copies/mL and  $\gamma=23$  day<sup>-1</sup>. Parameter  $\alpha_S$  was assumed identical for all individuals with  
 656 estimated value 42.7 cells day<sup>-1</sup>  $\mu L^{-1}$ .

id	t_det	$\delta_S$	log <sub>10</sub> $\beta$	log <sub>10</sub> $\pi$	$\kappa$	$\eta$
10066	9.62324	0.00015132	-4.65193	0.909672	0.277224	0.123357
10203	8.54516	0.00204632	-3.26718	0.817297	0.527654	0.127292
10220	6.65122	0.00313841	-3.08099	0.620871	0.401491	0.114859
10428	11.1077	0.00073429	-4.59846	1.54429	0.368804	0.111938
10435	19.0919	0.00027166	-4.47328	0.904101	0.264548	0.110447
10463	6.84935	0.00026507	-4.38649	0.889934	0.301887	0.131843
10723	9.97963	0.00146281	-3.99989	1.37432	0.498574	0.141124
10739	9.411	0.00019097	-4.95646	1.24387	0.261739	0.127887
10742	10.3869	0.00082387	-5.0462	2.02491	0.37825	0.1278
40007	6.24852	0.00018223	-5.21088	1.54045	0.294554	0.163842
40061	5.39621	0.00314832	-3.31239	0.903244	0.445166	0.141676
40094	12.6378	0.00199422	-4.13477	1.32797	0.309711	0.13599
40100	11.2924	0.00025118	-5.36688	1.82283	0.308901	0.140513
40123	7.57667	0.00329133	-3.84293	1.30979	0.365231	0.121762
40168	13.9136	0.00060717	-4.75154	1.43615	0.265229	0.168813
40231	9.42558	0.00076336	-4.20229	1.12247	0.335064	0.127668
40250	7.36099	0.00109501	-4.47884	1.52832	0.348449	0.147586
40257	6.78964	0.00011896	-5.1431	1.42538	0.343796	0.168909
40265	6.95792	0.00136279	-4.2837	1.42853	0.357415	0.156733
40353	6.31939	0.00143253	-3.2951	0.661891	0.47208	0.115288
40363	6.62144	0.00146629	-4.34715	1.55672	0.385113	0.139406
40511	5.98824	0.00059542	-4.18057	1.39129	0.633979	0.165386
40512	7.61647	0.00135145	-4.56775	1.71063	0.365675	0.151776
40577	9.63581	0.00013342	-5.38839	1.54198	0.263169	0.151545
10753	6.92639	0.00026048	-4.84951	1.30014	0.288365	0.179543

40032	9.56705	0.00444984	-4.01053	1.99893	0.763984	0.143641
40211	13.6817	0.00058183	-4.29235	1.307	0.462561	0.173981
40503	6.45906	0.00070623	-4.67312	1.57485	0.35149	0.156667
40646	5.3783	0.00013301	-4.82366	1.06741	0.294896	0.131254
10204	17.0122	0.00036347	-4.53213	1.08932	0.280075	0.116799
10374	10.5916	0.00273136	-3.58615	1.49338	0.847886	0.21121
40067	16.4039	0.00214128	-4.33294	1.95401	0.629711	0.1607
40096	11.6224	0.00035585	-4.74625	1.35842	0.316595	0.145131
40134	13.9472	0.00026626	-4.17889	0.794832	0.359614	0.143679
40139	22.4029	0.00150565	-3.88411	1.71941	1.07298	0.103824
40242	20.7467	0.00125992	-3.90151	1.03766	0.357604	0.126244
40283	13.8411	0.000963	-4.25141	1.62302	0.640396	0.106864
40195	19.2263	0.00180051	-4.28104	1.76886	0.550271	0.169485
40435	18.0127	0.00076233	-4.39397	1.53545	0.497526	0.190491
40492	12.2856	0.00071631	-4.55015	1.37054	0.300422	0.199821
40528	2.51471	0.011889	-3.0486	2.29511	3.53815	0.213035
40652	15.3836	0.00036841	-3.87475	0.703953	0.421935	0.173104
40700	24.9363	0.00383498	-4.09306	1.66294	0.407487	0.160303
40436	10.8677	0.00066651	-4.11294	0.839844	0.257823	0.173444
40491	13.3018	0.00023218	-4.43882	0.892541	0.299901	0.142118
40737	32.6352	0.00028113	-4.51066	0.978777	0.27742	0.128034

657

658 **Supplementary Table 3.** Best estimate of infection time relative to first positive viral load for each  
 659 individual. This table illustrates 5 previously reported methods<sup>3</sup> and our population nonlinear mixed  
 660 effects (pNLME) viral dynamics modeling approach. These methods are the maximum slope of any  
 661 two points on the upslope (max\_slope), the best log-linear regression slope (linear\_model), self-  
 662 reported entries from trial participants (diary), Bayesian phylogenetic inference of median time to  
 663 most-recent common ancestor (BEAST), and Poisson fitter diversity estimate assuming star-like  
 664 phylogeny.

id	max-slope	linear_model	Diary	BEAST	PFitter	pNLME
10066	-7.0240004	-18.176592		-2.62	-8	-16.372267
10203	-10.247329	-25.198679				-13.74697
10220	1.18944995	-2.6721271			-13	-16.04084
10428	-18.571441	-19.277159		-12.18	-17	-14.198643
10435	-27.053564	-92.332248				-26.6943
10463	-4.1810206	-7.290377		-25.42	-7	-12.977523
20225	-34.189189	-34.189189		-4.56	-6	
20245						
20263	-4.3846261	-7.4981795				
20314	-20.822782	-23.49317		-3.25	-10	
20337	-5.4536831	-11.143324				
20355	-10.400144	-22.516374				
20368	-3.07086	-7.8134674		1.87	-2	
20442	-5.2628089	-8.7491781				
20502	-7.7774033	-13.562055			-4	
20507	-5.4597332	-9.8952878		-9.64	-7	
20509	-8.4397215	-17.466691		-8.9	-6	
20511	-15.132734	-31.487589		-0.9	-2	
20631	-12.120706	-12.731752		-11.29	-30	

It is made available under a [CC-BY-ND 4.0 International license](#) .

<b>30112</b>	-3.4486448	-10.198418		-10.97	-18	
<b>30124</b>	-6.9575257	-15.019169			-16	
<b>30190</b>	-6.0162679	-19.733117		-1.35	-6	
<b>30507</b>	-13.070499	-21.902597				
<b>30812</b>	-5.3750022	-10.439993				
<b>30924</b>	-26.677419	-26.677419		-7.82	-11	
<b>40007</b>	-4.0959219	-8.7359897		-6.17	-4	-10.715117
<b>40061</b>	-1.3961514	-4.0645918	-13	0.39	-7	-13.25748
<b>40094</b>	-9.3396373	-28.128959		-19.12	-20	-22.087233
<b>40100</b>	-15.397661	-16.900742	-8	1.37	2	-15.072733
<b>40123</b>	-2.5469005	-5.6085575	-7		-36	-15.28629
<b>40168</b>	-0.3014733	-2.8255718	-4	-6.87	-6	-20.232067
<b>40231</b>	-9.2630352	-18.909093	-7	-6.45	-9	-14.589803
<b>40250</b>	-5.1308885	-10.832909	-15	-5.69	-2	-12.941387
<b>40257</b>	-3.714285	-6.8801494	-9	-6.1	-11	-9.51851
<b>40265</b>	-3.6988264	-6.4509726	-8	-21.74	-19	-13.166543
<b>40353</b>	-0.8379281	-3.7858275	-3	-0.46	6	-13.518597
<b>40363</b>	-3.4435622	-7.4401583	-9		-2	-11.435493
<b>40436</b>	-6.498836	-19.142592	-5	1.14	-6	-18.718233
<b>40491</b>	-14.437486	-28.648138	-14			-17.9225
<b>40511</b>			-3	-5.65	-18	-8.7821767
<b>40512</b>	-5.2047177	-11.910279	-20	2.62	2	-12.201273
<b>40577</b>	-7.5341659	-16.900867		-6.39	-14	-14.813207

665