



HAL
open science

Model Averaging in Viral Dynamic Models

Antonio Gonçalves, France Mentré, Annabelle Lemenuel-Diot, Jérémie Guedj

► **To cite this version:**

Antonio Gonçalves, France Mentré, Annabelle Lemenuel-Diot, Jérémie Guedj. Model Averaging in Viral Dynamic Models. *AAPS Journal*, 2020, 22 (2), pp.48. 10.1208/s12248-020-0426-7. inserm-02617421

HAL Id: inserm-02617421

<https://inserm.hal.science/inserm-02617421>

Submitted on 25 May 2020

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Model Averaging in Viral Dynamic Models

Antonio Gonçalves¹, France Mentré¹, Annabelle Lemenuel-Diot² and Jérémie Guedj¹

¹Université de Paris, IAME, INSERM, F-75018 Paris, France

²Roche Pharmaceutical Research and Early Development, Pharmaceutical Sciences, Roche
Innovation Center Basel

Corresponding author: Antonio Gonçalves, antonio.goncalves@inserm.fr

IAME INSERM U1137,

16 rue Henri Huchard 75018, Paris, France

Tel: +33 1 57 27 75 39

13 **Abstract**

14 The paucity of experimental data makes both inference and prediction particularly challenging
15 in viral dynamic models. In presence of several candidate models, a common strategy is
16 model selection (MS), in which models are fitted to the data but only results obtained with the
17 “best model” are presented. However, this approach ignores model uncertainty, which may
18 lead to inaccurate predictions. When several models provide a good fit to the data, another
19 approach is model averaging (MA) that weights the predictions of each model according to its
20 consistency to the data.

21 Here we evaluated by simulations in a nonlinear mixed-effect model framework the
22 performances of MS and MA in two realistic cases of acute viral infection: i) inference in
23 presence of poorly identifiable parameters, namely initial viral inoculum and eclipse phase
24 duration ii) uncertainty on the mechanisms of action of the immune response.

25 MS was associated in some scenarios with a large rate of false selection. This led to a
26 coverage rate lower than the nominal coverage rate of 0.95 in the majority of cases and below
27 0.50 in some scenarios. In contrast, MA provided better estimation of parameter uncertainty,
28 with coverage rates ranging from 0.72 to 0.98 and mostly comprised within the nominal
29 coverage rate. Finally, MA provided similar predictions than those obtained with MS.

30 In conclusion, parameter estimates obtained with MS should be taken with caution, especially
31 when several models well describe the data. In this situation, MA has better performances and
32 could be performed to account for model uncertainty.

33

34 **Introduction**

35 Since 1995 and the two seminal papers providing an estimate of the half-life of HIV
36 particles in blood (1,2), the use of viral dynamic models has considerably expanded.
37 Applications have been summarized in a recent issue of Immunological Reviews (3), showing
38 their relevance for understanding the host-pathogen interactions in both chronic and acute
39 infections (4–6). In the last decade, parameter estimation of these models has increasingly
40 relied on nonlinear mixed effect models (NLMEM), a statistical approach that improves both
41 precision and accuracy of estimates by explicitly taking into account the between-subjects
42 variability in the model (7,8). This is particularly true in the case of antiviral drug
43 development where NLMEM have become central to support optimal treatment strategies in
44 presence of a large variability in the response (9,10).

45 Although inference has been greatly facilitated by the use of NLMEM, viral dynamic models
46 remain often characterized by a lack of theoretical or practical identifiability (7,8). In fact the
47 availability of powerful algorithms for inference has mechanically led to the development of
48 increasingly complex models, questioning the reliability of viral kinetic parameter estimates.
49 In order to improve identifiability of these models, a commonly used strategy is to fix
50 parameters to plausible values and then to check the impact of these choices by conducting
51 sensitivity analyses. For instance in acute viral infection, one can fix the initial viral inoculum
52 or the eclipse phase duration, two parameters that can hardly be estimated using only viral
53 load data (11,12). Data fitting can also be used to evaluate the plausibility of different
54 biological assumptions. In that case the usual approach is model selection (MS), where a
55 predefined set of candidate models are fitted to the data and the model providing the best fit to
56 the data (based on Akaike or Bayesian Information Criteria) is selected and carried forward in
57 the analysis. In both contexts, these approaches, by focusing the predictions on a single

58 model, ignore the model uncertainty and may lead to wrong predictions (13) and potentially
59 inaccurate biological conclusions (14–16).

60 In this paper, we propose to use model averaging (MA) as an alternative approach to MS in
61 viral dynamics. MA is a conceptually simple approach, where the uncertainty related to each
62 candidate model is taken into account and predictions associated to each model are weighted
63 based on their consistency with the data (17,18). Through an extensive simulation study, we
64 compare parameter estimates and predictive performances of model averaging versus model
65 selection. We discuss the benefits and limits of model averaging compared to model selection.
66 Simulations are inspired from recent works in Zika and Ebola virus dynamics (19,20)
67 representing two typical settings encountered in viral dynamic modeling: i) a set of
68 parameters are fixed to arbitrary values of a given biological model to ensure identifiability;
69 ii) model selection relies on the comparison of fitting criterion of a set of pre-defined different
70 biological models.

71

72 **Material and methods**

73 **Model selection and model averaging**

74 **Model for the observations.** Let Y_{ijm} denote the j^{th} log viral load measurement of subject i at
75 time j , and suppose that $m=1,\dots,M$ candidate models can be used to simulate the data. The
76 model for the observations is defined as:

$$77 \quad Y_{ijm} = \log_{10}[V_m(t_j, \theta_{im})] + e_{ijm} \quad (1)$$

78 where V_m is the viral load prediction function given by model m , θ_{im} is the vector of
79 individual parameters under model m , t_j the time of viral load measurement, assumed to be
80 similar for all patients and all models, and e_{ijm} the residual error. Individual parameters θ_{im}
81 are log-normally distributed and depend on the vector of fixed effects μ_m and the vector of
82 random effects $\eta_{im} \sim \mathcal{N}(0, \Omega_m)$ with $\theta_{im} = \mu_m \times e^{\eta_{im}}$. The variance-covariance matrix Ω_m is
83 assumed to be diagonal. Residuals errors are assumed to be independent and normally
84 distributed $e_{ijm} \sim \mathcal{N}(0, \sigma_m^2)$. Each biological model m is therefore associated with a set of
85 population parameters, Ψ_m , of dimension $p_m = \dim(\mu_m, \Omega_m, \sigma_m)$.

86 **Inference and model selection.** For each candidate model, one can estimate the parameters
87 using maximum likelihood estimates, providing, for each model, an estimate of the population
88 parameters, noted $\hat{\Psi}_m$. One can also provide the confidence intervals of the parameters of
89 interest under each model. This can be done using the asymptotic approximation where the
90 density function of the estimated parameter $\hat{\Psi}_m$, $p(\hat{\Psi}_m)$, is assumed to be Gaussian with a
91 variance-covariance matrix given by the inverse of the Fisher Information Matrix (FIM^{-1}).

92 Then, the most common approach is to select the model that best describes the data. This can
93 be done using various criteria that rely on penalizations of the log-likelihood (LogL), such as

94 the Akaike information criteria (AIC), the consistent Akaike (CAIC) or the Bayesian
95 information criteria (BIC) (21–24). In line with previous analysis (25,26), we relied on AIC
96 given by $AIC_m = -2\text{LogL}(\hat{\Psi}_m) + 2p_m$. The analysis then focuses on the results (i.e., parameter
97 estimates, confidence intervals, and predictions) obtained with the “best” model, i.e., the
98 model associated with the lowest AIC among the m candidate models, noted AIC_{\min} , with
99 parameter estimates noted $\hat{\Psi}_{MS}$.

100 **Model averaging.** As explained above, MS is limited in the sense that it ignores the
101 uncertainty associated with each model and only focuses on a post hoc selected model (14).

102 Alternatively, one can use model averaging (MA) to take into account the fact that several
103 candidate models may provide a reasonable fit to the data . In this approach a weight is

104 attributed to each candidate model, w_m , proportional to AIC, such as $w_m = \frac{e^{-\frac{\Delta AIC_m}{2}}}{\sum_{m=1}^M e^{-\frac{\Delta AIC_m}{2}}}$ where

105 $\Delta AIC_m = AIC_m - AIC_{\min}$ (14,17,18). In that case the MA estimator of Ψ_m is given by $\hat{\Psi}_{MA}$,

106 with a density function given by $p(\hat{\Psi}_{MA}) = \sum_{m=1}^M w_m p(\hat{\Psi}_m)$. Another approach could be to

107 consider only the models that are responsible for the majority of the weight (0.9 or 0.8), and

108 equally average them (see Discussion section).

109

110 **Viral dynamic settings**

111 Our objective is to compare model selection and model averaging in two typical contexts of

112 viral dynamic models. In the first setting, we focus on the issue arising from using model

113 selection when some parameters of the model cannot be identified and are fixed to arbitrary

114 values. In the second setting, we focus on the issue arising from using model selection when

115 several different biological models can be proposed to fit the data.

116 **Setting I: viral dynamic models in presence of poorly identifiable parameters.** We here
117 focus on the standard target cell limited (TCL) model given by:

$$\frac{dT}{dt} = -\beta TV \quad (2)$$

$$\frac{dI_1}{dt} = \beta TV - kI_1 \quad (3)$$

$$\frac{dI_2}{dt} = kI_1 - \delta I_2 \quad (4)$$

$$\frac{dV}{dt} = \pi I_2 - cV \quad (5)$$

$$T_{t=0} = T_0; I_{1,t=0} = 0; I_{2,t=0} = 0; V_{t=0} = V_0$$

118 where, T are the target cells, I_1 the infected cells in eclipse phase, I_2 the productive infected
119 cells and V the viral load in plasma. The model depends on the following disease parameters:
120 β the infectivity rate constant, k the eclipse rate, δ the infected cell elimination rate, π the viral
121 production rate constant, c the clearance of free virus, T_0 the initial number of target cells and
122 V_0 the initial viral load. For the ease of interpretation and fitting, we reparametrized the model
123 as $R_0 = \frac{\beta\pi T_0}{c\delta}$, the basic reproductive ratio, instead of β , where R_0 represents the number of
124 secondary infection caused by one infected cells when the target cells are abundant. For the
125 sake of simplificty we focused here on a simple, exponentially distributed, duration for the
126 eclipse phase, but more complex models can be considered (27).

127 Not all parameters of the TCL model can be uniquely identified when only the viral load data
128 are available (28–30) and this issue is not circumvented when parameters are estimated using
129 NLMEM. This can be shown by analyzing the expected standard errors obtained with the
130 approximated Fisher Information Matrix (<http://www.pfim.biostat.fr/>; see more details in
131 (30)). Table I provides the expected standard errors obtained with 30 individuals sampled 3
132 days from day 3 up to day 18 post infection using typical parameter values close to those
133 found during Zika infection in nonhuman primates (19). Although being theoretically
134 identifiable, several parameters are associated with a very large expected standard error. This

135 can be corrected by fixing V_0 and k to some arbitrary values, leading to expected relative
136 standard errors lower than 30% for all parameter values. We here aim to evaluate the impact
137 of the choice of k and V_0 on parameter estimates.

138 We also conducted a sensitivity analysis on the parameters variability of setting I ($k=4 \text{ d}^{-1}$ and
139 $V_0=10^{-4} \text{ copies.mL}^{-1}$). We reported the influence of lower and higher variabilities ($\omega=0.1$ and
140 $\omega=1$, respectively) on the performances of MS and MA.

141

142 **Setting II: viral dynamic models including the immune response.** In order to evaluate the
143 impact of testing different biological assumptions in parameter estimates and predictions, we
144 considered 4 additional models integrating the role of innate or adaptive immune response in
145 the control of viral replication and inspired from the models used to describe Ebola infection
146 in nonhuman primates (20). These models extend the TCL model with an additional
147 compartment, noted F. This compartment is not observed and can therefore represent any
148 biological entity involved in viral clearance, such as cytokine, macrophages, T-cell or
149 antibodies (5,32,33). We assumed that F is produced at a rate q proportional to the number of
150 productively infected cells, I_2 , and is eliminated at a rate d_F (20). Thus F could either i)
151 increase the number of refractory infected cells (refractory model, R), ii) decrease the viral
152 production (production inhibition model, PI), iii) increase the clearance of productive infected
153 cells (cytotoxic model, C) or iv) increase the clearance of the virus (virus-killing model, V).
154 In all models, the effects of F followed an E_{\max} relationship with ϕ the maximal effect of F
155 and θ the sensitivity parameter. Table II displays the four model equations.

156

157 **Simulations and parameter estimation**

158 For each setting, the simulation procedure and parameter estimation under both MS
159 and MA are described below.

160 **Parameter values.** In the first setting, we aimed to evaluate the impact of fixing the two
161 poorly identifiable parameters, k and V_0 , in the target cell limited model. For that purpose we
162 defined a set of $M=9$ candidate models with values for V_0 and k equal to $V_0 = 10^{-5}; 10^{-4}$ or 10^{-3}
163 copies.mL^{-1} and $k = 1; 4$ or 20 d^{-1} . The other parameter values are given in Table III. Figure
164 1A shows that viral load is biphasic with a peak close to $10 \log_{10} \text{copies.mL}^{-1}$ in all 9
165 scenarios, but the time to peak depends on k and V_0 to a lesser extent.

166 In the second setting V_0 and k were fixed to $10^{-4} \text{ copies mL}^{-1}$ and 4 d^{-1} , respectively. To
167 ensure a fair comparison between the models, parameters were chosen to predict a similar
168 contribution of the immune response to viral control, as measured by the area under the curve
169 of the \log_{10} viremia from 0 to 20 days (AUC). Thus, in all four models, the parameter values
170 were such that $\text{AUC}=100 \log_{10} \text{copies.days.mL}^{-1}$ while assuming $\phi^* = 0$ would lead to AUC
171 $120 \log_{10} \text{copies.days.mL}^{-1}$ (i.e., the absence of an immune system would lead to a 20%
172 increase in AUC) (Figure 1B). The values of the TCL model were chosen to lead to a similar
173 $\text{AUC}=100 \log_{10} \text{copies.days.mL}^{-1}$.

174 We assumed that R_0 , δ , π , θ and ϕ were the estimated fixed effects. Those parameters,
175 with the exception of θ , were associated to an intermediate between-subject variability, ω ,
176 equal to 0.3. Other parameters were assumed to be known with values given in Table III.

177 **Data simulation.** For each model we simulated $S=300$ datasets of $N=30$ individuals using the
178 same population parameter values, Ψ_m^* , given in Table III. Therefore $S \times M = 2700$ datasets
179 were simulated in the setting I, and 1500 datasets were simulated in the setting II. We
180 assumed measurements were made at days 3, 6, 9, 12, 15 and 18, with a lower limit of

181 quantification (LLOQ) of $1 \log_{10} \text{copies.mL}^{-1}$ (20), and a measurement error term, σ , equal to
182 $0.7 \log_{10} \text{copies.mL}^{-1}$.

183 **Parameter estimation.** Each of the $s=1, \dots, S$ dataset, was fitted using the M candidate
184 models of each setting. The set of parameter estimates obtained on the dataset s using the
185 model m , namely R_0 , δ , π , θ and ϕ and their corresponding between-subjects variabilities if
186 specified, was noted without loss of generality $\hat{\Psi}_m^s$. Parameter estimates were obtained by
187 maximization of the likelihood using the SAEM algorithm implemented in the MONOLIX
188 software (version 2018, release 2). We used $k_1=800$ and $k_2=200$ iterations for the exploratory
189 and smoothing phases, respectively. We used the asymptotic approximation to derive the
190 probability density function of $\hat{\Psi}_m^s$, noted $p(\hat{\Psi}_m^s)$, assumed to be Gaussian with a variance-
191 covariance matrix given by the inverse of the Fisher Information Matrix (FIM^{-1}). The FIM
192 was computed by stochastic approximation with at least 100 and up to 800 iterations. Of note,
193 among the M models used to fit the data, only one is the true model (noted TM), i.e., the
194 model used to generate the data, and we note $\hat{\Psi}_{TM}^s$ the parameter estimates obtained by fitting
195 the dataset s with TM.

196 For MA, 95% confidence intervals of $\hat{\Psi}_{MA}^s$ was then calculated by sampling 10,000 values in
197 the mixture distribution $p(\hat{\Psi}_{MA}^s) = \sum_{m=1}^M w_m^s p(\hat{\Psi}_m^s)$ and computing the associated 2.5th and
198 97.5th percentiles (14,18,34).

199

200 **Performances of model averaging and model selection for estimation**

201 **Model selection.** For each scenario, we reported the distribution weight of each candidate
202 model as well as the proportion of simulations where each candidate model was selected
203 (based on AIC, see above).

204 **Parameter estimates and comparison with true parameter value.** For each scenario, we
205 reported the coverage rate obtained for each parameter with estimator based on MS, MA or
206 TM, defined as the proportion of simulated datasets for which the true value of the parameter
207 was contained in the 95% confidence interval of the estimated parameter. The coverage rates
208 were compared with the prediction interval of a Binomial distribution with $p=0.95$ and $S=300$,
209 i.e., $[0.923; 0.973]$ and were reported for parameters R_0 , π and δ in setting I and R_0 and δ in
210 setting II.

211

212 **Performances of model averaging and model selection for prediction**

213 Finally, we aimed to evaluate MA in the context of prediction, i.e., the capability to anticipate
214 the effect of a change in the experimental setting. We focused on the prediction of the impact
215 of an antiviral treatment limiting the viral production π with efficacy ε ($0 < \varepsilon < 1$) on the
216 predicted proportion of patients with undetectable viral load ($10 \text{ copies mL}^{-1}$) at a given time
217 point. We assumed that treatment was initiated at time $t=6$ and lasted until $t=20$ days, which
218 coincides with the end of the follow-up. We considered 3 levels of efficacy on decreasing
219 viral production with a factor $1 - \varepsilon$, namely $\varepsilon = 0.90, 0.95$ and 0.99 , and we focused on the
220 prediction at $t=20$. For each model and each value of ε , Monte-Carlo simulations were used to
221 the expected proportion of patients below the limit of detection noted $P_m^{*,\varepsilon}(\%) =$
222 $P[V_m(t = 20, \Psi_m^*, \varepsilon) < 10]$.

223 Following what has been done above, one can calculate, for each simulated dataset, the
224 estimate that would be given by model selection, given by $P_{MS}^{S,\varepsilon} = P[V_{MS}(t = 20, \hat{\Psi}_{MS}^S, \varepsilon) <$
225 $10]$ or by model averaging $P_{MA}^{S,\varepsilon} = \sum_{m=1}^M w_m^S P[V_m(t = 20, \hat{\Psi}_m^S, \varepsilon) < 10]$. Likewise for the
226 sake of comparison, one can also calculate the probability obtained by fitting the data under

227 the true model, $P_{TM}^{s,\varepsilon} = P[V_m(t = 20, \hat{\Psi}_{TM}^s, \varepsilon) < 10]$. These values were summarized by
228 calculating the bias and root mean square error (RMSE), given by $\frac{1}{S}\sum(P_{MA}^{s,\varepsilon} - P_m^{*,\varepsilon})$ and
229 $\sqrt{\frac{1}{S}\sum(P_{MA}^{s,\varepsilon} - P_m^{*,\varepsilon})^2}$ in the case of MA (similar applies to calculate the bias and RMSE in the
230 case of MS or TM). For the sake of graphical representation, proportions of patients below the
231 limit of detection were presented as percentages and biases and RMSE were therefore
232 expressed in percentages.

233

234 **Results**

235 **Setting I**

236 The first setting focused on the comparison between model averaging and model
237 selection when parameters of the model (e.g., the eclipse phase, k , and the initial inoculum V_0)
238 cannot be identified and are fixed to arbitrary values (see Table III).

239 Overall, the true set of parameter values was selected up to 62% of the simulations. The
240 two parameters did not have the same rate of selection, with the correct values for k and V_0
241 being selected up to 71% and 96%, respectively (Figure 2A). Although the true model was not
242 systematically associated with the lowest AIC, it was associated in all scenarios with the
243 largest weight among the candidate models with a median value comprised between 0.32 and
244 0.55 (Figure 2B). In all cases considered, at least two models had a weight greater than 0.20.

245 We next evaluated the impact of these results on parameter estimates and coverage rates.
246 The estimation of R_0 using model selection was associated with a poor coverage rate between
247 0.46 and 0.63 (Figure 3). Results for the loss rate of infected cells, δ , were better with a
248 coverage rate ranging from 0.53 and 0.94, and was comprised in the nominal 0.95 coverage

249 rate in 3 out of 9 scenarios. For the viral production π , MS showed coverage rates ranging
250 from 0.68 and 0.96, and was comprised in the nominal 0.95 coverage rate in 5 out of 9
251 scenarios. Model averaging largely improved the coverage rates for all parameters and gave
252 results close to those obtained with the true model. The coverage rates was between 0.91 and
253 0.98 for R_0 , between 0.72 and 0.95 for δ and between 0.78 and 0.98 for π . Further the
254 coverage rates were comprised in the nominal 0.95 coverage rate in 7 out 9 scenarios for R_0
255 and 5 out of 9 scenarios for δ and π . All confidence intervals can be found in supplemental
256 figures S1, S2 and S3.

257 Lastly we explored the effect of simulating with less ($\omega=0.1$) or more ($\omega=1$) inter-
258 subjects variability. In both cases, MS provided subnominal coverage rates but MA corrected
259 them (see Supplemental Figure S4). Eventually, we observed poorer coverage rates with
260 $\omega=0.1$ and improved with $\omega=1$ compared to $\omega=0.3$.

261

262 **Setting II**

263 In the second setting, we assessed the properties of parameter estimates when several
264 biological models can be proposed. We focused on models characterizing the effect of the
265 immune response, considering that the immune response compartment could alternatively
266 make cells refractory to infection, limit the production of virus, increase the elimination of
267 infected cells or increase the elimination of free virions (see Table III).

268 Unlike what was found in the previous setting, the chance of selecting the true model was
269 largely dependent on the model considered. In fact these chances were equal to 97% for the
270 refractory model but this percentage could decrease to 58% with the cytotoxic model (Figure
271 4A). Conversely, the models were also associated with a large rate of false selection with rates

272 ranging from 3 to 10% for the refractory model, and up to 20% for the production inhibition
273 model. In the case of the target cell limited model, the chances of correctly selecting it were
274 equal to 88% and the rate of false selection were ranging from 1 to 19%. The median weight
275 associated to the true model ranged from 0.43 to 0.99 (Figure 4B).

276

277 Accordingly, MS provided satisfactory coverage rates for target cell limited and
278 refractory model (Figure 5); however it failed to achieve the nominal coverage rate in all other
279 models, with values ranging from 0.62 to 0.91 for R_0 and from 0.50 to 0.89 for δ . This could
280 be improved by taking into account model uncertainty and using model averaging. Indeed the
281 coverage rates ranged from 0.86 to 0.99 for both parameters in all models considered. In fact,
282 MA had even better performances than the true model in some cases, which achieved
283 subnominal coverage rate in 3 of the 5 considered scenarios (Figure 5). All confidence
284 intervals can be found in supplemental figures S5 and S6.

285

286 Finally, we compared the predictive performances of MS, MA and the true model. For
287 that purpose, we predicted the effect of a putative antiviral treatment on the proportion of
288 patients having undetectable viremia at end of follow up (day 20). Here as well the
289 performances obtained using model selection and model averaging were compared. In all
290 cases, the percentage of undetectable viral loads at end of treatment was accurately predicted
291 for both MS and MA, with no more than 4% of bias in all cases considered (Figure 6). In term
292 of precision of estimation, the results were also largely similar in most scenarios, with RMSE
293 ranging from 0.4 to 30.5% in all cases. In one case, namely $\varepsilon=0.95$, we found that MA
294 outperformed the results obtained by MS. Here as well, the results obtained by model
295 averaging were largely comparable with those obtained with the true model.

296 **Discussion**

297 The objective of this study was to compare the estimation and the predictive
298 performances of model selection and model averaging in the context of viral dynamic models.
299 We explored two frequent issues encountered when developing viral dynamic models with
300 uncertainty related either to (I) unidentifiable parameters or (II) the presence of several
301 candidate biological models. In the two settings MS provided poor coverage rates of typical
302 parameters. This stems from the fact that MS neglects model uncertainty and focuses on one
303 single “best model”, leading to overconfidence in the parameter estimates. This can be
304 corrected under certain conditions using MA, which provided better coverage rates and
305 achieved the nominal coverage rate in most scenarios studied. MA can also be relevant to
306 predict the effect of intervention, such as the percentage of patients that would achieve
307 undetectable viral loads during treatment. Thus extending results found in other contexts, in
308 particular dose finding studies (25,35).

309 By offering a simple framework to take into account model uncertainty, MA accounts for the
310 fact that in many situations several biological models are plausible. Our study shows the
311 limitations of reporting only the best model. For instance, in the case of the target cell limited
312 model, we found that the chance to conclude wrongly to an immune response controlling the
313 infection was equal to 11%. In the case of the refractory model, which has been proposed as a
314 driving force in several acute infection (33,36), our results were more reassuring, with a rate
315 of false rejection of only 3%. This risk was larger with other models integrating an immune
316 response, with rates of false rejection greater than 60% in some cases. By weighing the
317 predictions of alternative models, MA avoids the caveat of MS. As advocated in other
318 contexts (37,38), MA can be used to more transparently discuss model uncertainty and to
319 stimulate new data acquisition (13).

320 Although MA offers a simple alternative to MS, it also presents the defects of its
321 virtue. As MA weighs the models according to their information criterion, using MA is
322 relevant only if one model does not largely outperform the others. The weight value leading to
323 “outperformance” is arbitrary, and depends also on the number of candidate models.
324 Accordingly, one may question the need to use weights when making predictions. As
325 suggested by a reviewer, we conducted a simulation where all models having a weight greater
326 than a given threshold (0.1 or 0.2) were considered as equally likely in the prediction, and this
327 approach provided results close to those obtained with MA (Supplemental Figure S7).
328 MA still requires to make important assumptions that need to be kept in mind. First we used
329 the asymptotic Gaussian approximations to calculate the standard errors. This assumption
330 may not hold for all models, depending on their complexity and data paucity, as can be seen
331 in some cases of Figures 3 & 5. Other approaches have been proposed in the context of
332 NLMEM to calculate the standard error more precisely, such as bootstrapping (38), sampling
333 importance resampling (39) or Hamiltonian Monte-Carlo methods (HMC) (40). Future work
334 will be needed to evaluate in which contexts these methods, which are computationally
335 demanding, are beneficial. Second in our simulations, we assumed that there was a true model
336 and that it was part of the candidate models. Although there is no “true model” in real data,
337 we made this hypothesis to stress that MA should be performed only with biologically
338 relevant models. Likewise, MA should not be used to “blindly” average predictions of any
339 models and modelers should, prior to the analysis, develop other models at hand and, if
340 possible, discuss and perform new experiments to discriminate between them (13). In that
341 perspective, using MA to calculate CI is meaningful only if parameter have the same
342 interpretation across the candidate models. This is the case for half-life or viral production
343 rates but is less evident for derived parameters such the basic reproductive number R_0 (39,40).
344 Finally, MA does not substitute to a proper analysis of parameter identifiability. In fact the

345 differences between MA and MS may simply reveal a poor practical identifiability, i.e., the
346 fact that data available are not sufficient to precisely estimate parameters (8) and/or that the
347 biological question is wrongly formulated (13). This is also what we observed here, with the
348 wrong selection of models being in part due to the fact that the models had a poor practical
349 identifiability, at least for some parameters. In order to be performant, MA requires that only
350 a limited number of models are tested. It is only when a reasonable number of models remain
351 that MA can be relevant, as an alternative to Bayesian approaches, that may be tedious in
352 particular a non-linear mixed effect framework.

353

354 **Acknowledgments**

355 Antonio Gonçalves was funded by a grant from Roche Pharmaceutical Research and Early
356 Development. The authors declare that there is no conflict of interest regarding the
357 publication of this article. The authors also would like to acknowledge Hervé Le Nagard and
358 Lionel de la Tribouille for the use of CATIBioMed calculus facility.

359

360 **Supplemental information**

361 A readily usable R code to compute weights and confidence intervals in MA is provided in
362 supplementary material. The example is based on digitized Zika data of Best et al. *PNAS* 2017
363 and illustrates model averaging in the context of poorly identifiable parameters (setting I).

364

365 **References**

- 366 1. Ho DD, Neumann AU, Perelson AS, Chen W, Leonard JM, Markowitz M. Rapid
367 turnover of plasma virions and CD4 lymphocytes in HIV-1 infection. *Nature*. 1995 Jan
368 12;373(6510):123–6.
- 369 2. Wei X, Ghosh SK, Taylor ME, Johnson VA, Emini EA, Deutsch P, et al. Viral dynamics
370 in human immunodeficiency virus type 1 infection. *Nature*. 1995 Jan 12;373(6510):117–
371 22.
- 372 3. Perelson AS, Ribeiro RM. Introduction to modeling viral infections and immunity.
373 *Immunological Reviews*. 2018;285(1):5–8.
- 374 4. Perelson AS. Modelling viral and immune system dynamics. *Nat Rev Immunol*. 2002
375 Jan;2(1):28–36.
- 376 5. Best K, Perelson AS. Mathematical modeling of within-host Zika virus dynamics.
377 *Immunological Reviews*. 2018;285(1):81–96.
- 378 6. Ciupe SM. Modeling the dynamics of hepatitis B infection, immunity, and drug therapy.
379 *Immunological Reviews*. 2018;285(1):38–54.
- 380 7. Lavielle M, Mentré F. Estimation of population pharmacokinetic parameters of
381 saquinavir in HIV patients with the MONOLIX software. *J Pharmacokinet
382 Pharmacodyn*. 2007 Apr;34(2):229–49.
- 383 8. Guedj J, Thiébaud R, Commenges D. Practical identifiability of HIV dynamics models.
384 *Bulletin of Mathematical Biology*. 2007 Oct 25;69(8):2493–513.
- 385 9. Snoeck E, Chanu P, Lavielle M, Jacqmin P, Jonsson EN, Jorga K, et al. A
386 comprehensive Hepatitis C viral kinetic model explaining cure. *Clinical Pharmacology
387 & Therapeutics*. 2010 Jun;87(6):706–13.
- 388 10. Nguyen T, Guedj J. HCV kinetic models and their implications in drug development:
389 HCV kinetic models and their implications. *CPT: Pharmacometrics & Systems
390 Pharmacology*. 2015 Apr;4(4):231–42.
- 391 11. Handel A, Longini IM, Antia R. Towards a quantitative understanding of the within-host
392 dynamics of influenza A infections. *Journal of The Royal Society Interface*. 2010 Jan
393 6;7(42):35–47.
- 394 12. Smith AM, Adler FR, Ribeiro RM, Gutenkunst RN, McAuley JL, McCullers JA, et al.
395 Kinetics of coinfection with influenza A virus and streptococcus pneumoniae. Grenfell
396 BT, editor. *PLoS Pathogens*. 2013 Mar 21;9(3):e1003238.
- 397 13. Ganusov VV. Strong inference in mathematical modeling: a method for robust science in
398 the twenty-first century. *Frontiers in Microbiology* [Internet]. 2016 Jul 22 [cited 2019
399 Apr 4];7. Available from:
400 <http://journal.frontiersin.org/Article/10.3389/fmicb.2016.01131/abstract>

- 401 14. Buckland ST, Burnham KP, Augustin NH. Model selection: an integral part of inference.
402 Biometrics. 1997 Jun;53(2):603.
- 403 15. Boulesteix A-L. Ten simple rules for reducing overoptimistic reporting in
404 methodological computational research. Lewitter F, editor. PLOS Computational
405 Biology. 2015 Apr 23;11(4):e1004191.
- 406 16. Kirk PDW, Babbie AC, Stumpf MPH. Systems biology (un)certainties. Science. 2015
407 Oct 23;350(6259):386–8.
- 408 17. Burnham KP, Anderson DR. Model selection and multimodel inference: a practical
409 information-theoretic approach. 2. ed., [4. printing]. New York, NY: Springer; 2010. 488
410 p.
- 411 18. Claeskens G, Hjort NL. Model selection and model averaging [Internet]. Cambridge:
412 Cambridge University Press; 2008 [cited 2019 Sep 30]. Available from:
413 <http://ebooks.cambridge.org/ref/id/CBO9780511790485>
- 414 19. Best K, Guedj J, Madelain V, de Lamballerie X, Lim S-Y, Osuna CE, et al. Zika plasma
415 viral dynamics in nonhuman primates provides insights into early infection and antiviral
416 strategies. Proceedings of the National Academy of Sciences. 2017 Aug
417 15;114(33):8847–52.
- 418 20. Madelain V, Baize S, Jacquot F, Reynard S, Fizet A, Barron S, et al. Ebola viral
419 dynamics in nonhuman primates provides insights into virus immuno-pathogenesis and
420 antiviral strategies. Nature Communications. 2018 Dec;9(1).
- 421 21. Bertrand J, Comets E, Mentré F. Comparison of model-based tests and selection
422 strategies to detect genetic polymorphisms influencing pharmacokinetic parameters.
423 Journal of Biopharmaceutical Statistics. 2008 Nov 7;18(6):1084–102.
- 424 22. Bozdogan H. Model selection and Akaike's Information Criterion (AIC): The general
425 theory and its analytical extensions. Psychometrika. 1987 Sep;52(3):345–70.
- 426 23. Anderson DR, Burnham KP. Understanding information criteria for selection among
427 capture-recapture or ring recovery models. Bird Study. 1999 Jan;46(sup1):S14–21.
- 428 24. Neath AA, Cavanaugh JE. The Bayesian information criterion: background, derivation,
429 and applications: The Bayesian information criterion. WIREs Comp Stat. 2012
430 Mar;4(2):199–203.
- 431 25. Buatois S, Ueckert S, Frey N, Retout S, Mentré F. Comparison of model averaging and
432 model selection in dose finding trials analyzed by nonlinear mixed effect models. AAPS
433 J. 2018 29;20(3):56.
- 434 26. Aoki Y, Röshammar D, Hamrén B, Hooker AC. Model selection and averaging of
435 nonlinear mixed-effect models for robust phase III dose selection. J Pharmacokinet
436 Pharmacodyn. 2017 Dec;44(6):581–97.
- 437 27. Kakizoe Y, Nakaoka S, Beauchemin CAA, Morita S, Mori H, Igarashi T, et al. A
438 method to determine the duration of the eclipse phase for in vitro infection with a highly
439 pathogenic SHIV strain. Sci Rep. 2015 Sep;5(1):10371.

- 440 28. Xia X, Moog CH. Identifiability of nonlinear systems with application to HIV/AIDS
441 models. *IEEE Transactions on Automatic Control*. 2003 Feb;48(2):330–6.
- 442 29. Wu H, Zhu H, Miao H, Perelson AS. Parameter identifiability and estimation of
443 HIV/AIDS dynamic models. *Bulletin of Mathematical Biology*. 2008 Apr;70(3):785–99.
- 444 30. Miao H, Dykes C, Demeter LM, Cavanaugh J, Park SY, Perelson AS, et al. Modeling
445 and estimation of kinetic parameters and replicative fitness of HIV-1 from flow-
446 cytometry-based growth competition experiments. *Bulletin of Mathematical Biology*.
447 2008 Aug;70(6):1749–71.
- 448 31. Dumont C, Lestini G, Le Nagard H, Mentré F, Comets E, Nguyen TT, et al. PFIM 4.0,
449 an extended R program for design evaluation and optimization in nonlinear mixed-effect
450 models. *Computer Methods and Programs in Biomedicine*. 2018 Mar;156:217–29.
- 451 32. Baccam P, Beauchemin C, Macken CA, Hayden FG, Perelson AS. Kinetics of influenza
452 A virus infection in humans. *J Virol*. 2006 Aug;80(15):7590–9.
- 453 33. Pawelek KA, Huynh GT, Quinlivan M, Cullinane A, Rong L, Perelson AS. Modeling
454 within-host dynamics of influenza virus infection including immune responses. *PLoS*
455 *Comput Biol* [Internet]. 2012 Jun 28 [cited 2019 Jun 26];8(6). Available from:
456 <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3386161/>
- 457 34. Pinheiro J, Bornkamp B, Glimm E, Bretz F. Model-based dose finding under model
458 uncertainty using general parametric models. *Statist Med*. 2014 May 10;33(10):1646–61.
- 459 35. Schorning K, Bornkamp B, Bretz F, Dette H. Model selection versus model averaging in
460 dose finding studies: K. SCHORNING *ET AL*. *Statistics in Medicine*. 2016 Sep
461 30;35(22):4021–40.
- 462 36. Saenz RA, Quinlivan M, Elton D, MacRae S, Blunden AS, Mumford JA, et al.
463 Dynamics of influenza virus infection and pathology. *J Virol*. 2010 Apr;84(8):3974–83.
- 464 37. Hoeting JA, Adrian E. Raftery, Madigan D. Bayesian model averaging: a tutorial. *Statist*
465 *Sci*. 1999 Nov;14(4):382–417.
- 466 38. Burnham KP, Anderson DR. Model selection and multimodel inference: a practical
467 information-theoretic approach. 2. ed., [4. printing]. New York, NY: Springer; 2010. 488
468 p.
- 469 39. Lloyd AL. The dependence of viral parameter estimates on the assumed viral life cycle:
470 limitations of studies of viral load data. *Proceedings: Biological Sciences*.
471 2001;268(1469):847–54.
- 472 40. Ribeiro RM, Qin L, Chavez LL, Li D, Self SG, Perelson AS. Estimation of the initial
473 viral growth rate and basic reproductive number during acute HIV-1 infection. *Journal of*
474 *Virology*. 2010 Jun 15;84(12):6096–102.

475

476

477 **List of Tables**

478 **Table I: Expected standard error (SE) of the fixed effect parameters using a target cell**
 479 **limited model when the estimated parameters include or do not include the initial**
 480 **inoculum, V_0 , and the eclipse rate, k .** Expected standard errors were calculated using PFIM
 481 software and for a study design including 30 subjects sampled every 3 days from day 3 to day
 482 18.

Parameter (units)	Estimation of R_0 , δ , V_0 , k and π			Estimation restricted to R_0 , δ and π		
	Estimate	SE	Relative SE(%)	Estimate	SE	Relative SE(%)
R_0	12	62.0	516%	12	0.84	7.0%
δ (d^{-1})	1	0.10	10%	1	0.063	6.3%
π (copies.cell $^{-1}$.d $^{-1}$)	6000	3625	604%	6000	1446	24.1%
V_0 (copies.mL $^{-1}$)	10^{-4}	29.7	743%	10^{-4} (fixed)	-	-
k (d^{-1})	4	38.9	971%	4 (fixed)	-	-
c (d^{-1})	20 (fixed)	-	-	20 (fixed)	-	-
T_0 (cells.mL $^{-1}$)	10^8 (fixed)	-	-	10^8 (fixed)	-	-

483

484

485 **Table II: Differential equations system of immune response models.** At $t=0$ we have

486 $T_{t=0} = T_0; I_{1,t=0} = 0; I_{2,t=0} = 0$ and $V_{t=0} = V_0$.

487

	Target cell limited	Refractory	Production inhibition	Cytotoxic	Virus-killing
$\frac{dT}{dt} =$	$-\beta TV$	$-\beta TV - \frac{\phi TF}{F + \theta}$	$-\beta TV$	$-\beta TV$	$-\beta TV$
$\frac{dI_1}{dt} =$	$\beta TV - kI_1$	$\beta TV - kI_1$	$\beta TV - kI_1$	$\beta TV - kI_1$	$\beta TV - kI_1$
$\frac{dI_2}{dt} =$	$kI_1 - \delta I_2$	$kI_1 - \delta I_2$	$kI_1 - \delta I_2$	$kI_1 - \delta I_2 - \frac{\phi I_2 F}{F + \theta}$	$kI_1 - \delta I_2$
$\frac{dV}{dt} =$	$\pi I_2 - cV$	$\pi I_2 - cV$	$\pi \left(1 - \frac{\phi F}{F + \theta}\right) I_2 - cV$	$\pi I_2 - cV$	$\pi I_2 - cV - \frac{\phi F V}{F + \theta}$
$\frac{dF}{dt} =$	$qI_2 - d_F F$	$qI_2 - d_F F$	$qI_2 - d_F F$	$qI_2 - d_F F$	$qI_2 - d_F F$

488

489

490 **Table III: Parameter values used for simulations**

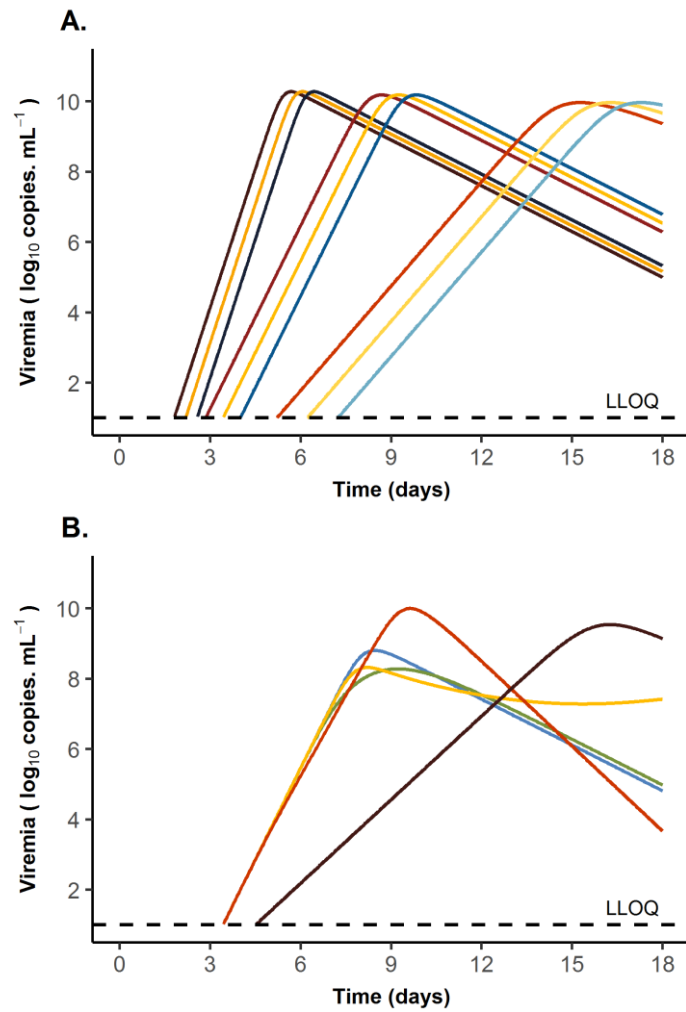
Parameter (units)	Setting I	Setting II				
	Target cell limited	Target cell limited	Refractory	Production Inhibition	Cytotoxic	Virus-killing
R_0^a	12	12				
δ^a (d ⁻¹)	1	1				
π^a (copie.cell ⁻¹ .d ⁻¹)	6000	250	6000			
ϕ^a	-	-	1	0.99	0.9	36.5
θ	-	-	2200	325000	3	0.001
V_0 (copies.mL ⁻¹)	{10 ⁻⁵ ; 10 ⁻⁴ ; 10 ⁻³ }	10 ⁻⁴				
k (d ⁻¹)	{1; 4; 20}	4				
c (d ⁻¹)	20	20				
T_0 (cells.mL ⁻¹)	10 ⁸	10 ⁸				
q (d ⁻¹)	1	1				
d_F (d ⁻¹)	0.4	0.4				
SD of the additive error	0.7	0.7				

^a: parameters for which inter-individual variability $\omega=0.3$

491

492

493 **Legend to Figures**



494

495 **Figure 1. Viral kinetics profiles obtained with the population parameters for each**

496 **candidate model. (A) and (B) correspond to the simulation settings I and II, respectively. In**

497 Panel A, curves are regrouped by 3. At left, the first 3 curves correspond to models with $k =$

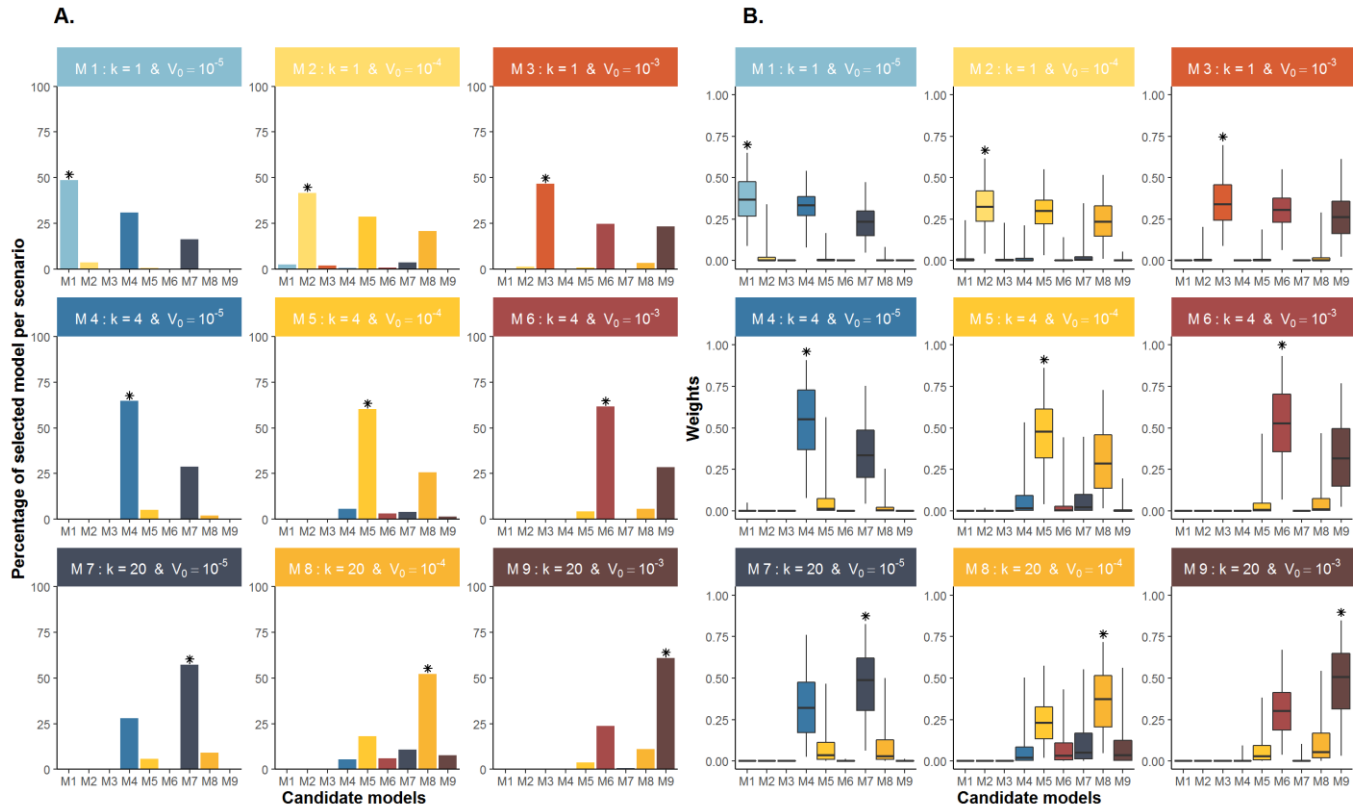
498 20 d^{-1} ; center, models with $k = 4 \text{ d}^{-1}$ and right, models with $k = 1 \text{ d}^{-1}$. Within each group, red

499 curves correspond to models with $V_0 = 10^{-5} \text{ copies.mL}^{-1}$, yellow curves to models with $V_0 =$

500 $10^{-5} \text{ copies.mL}^{-1}$ and blue curves to models with $V_0 = 10^{-5} \text{ copies.mL}^{-1}$.

501

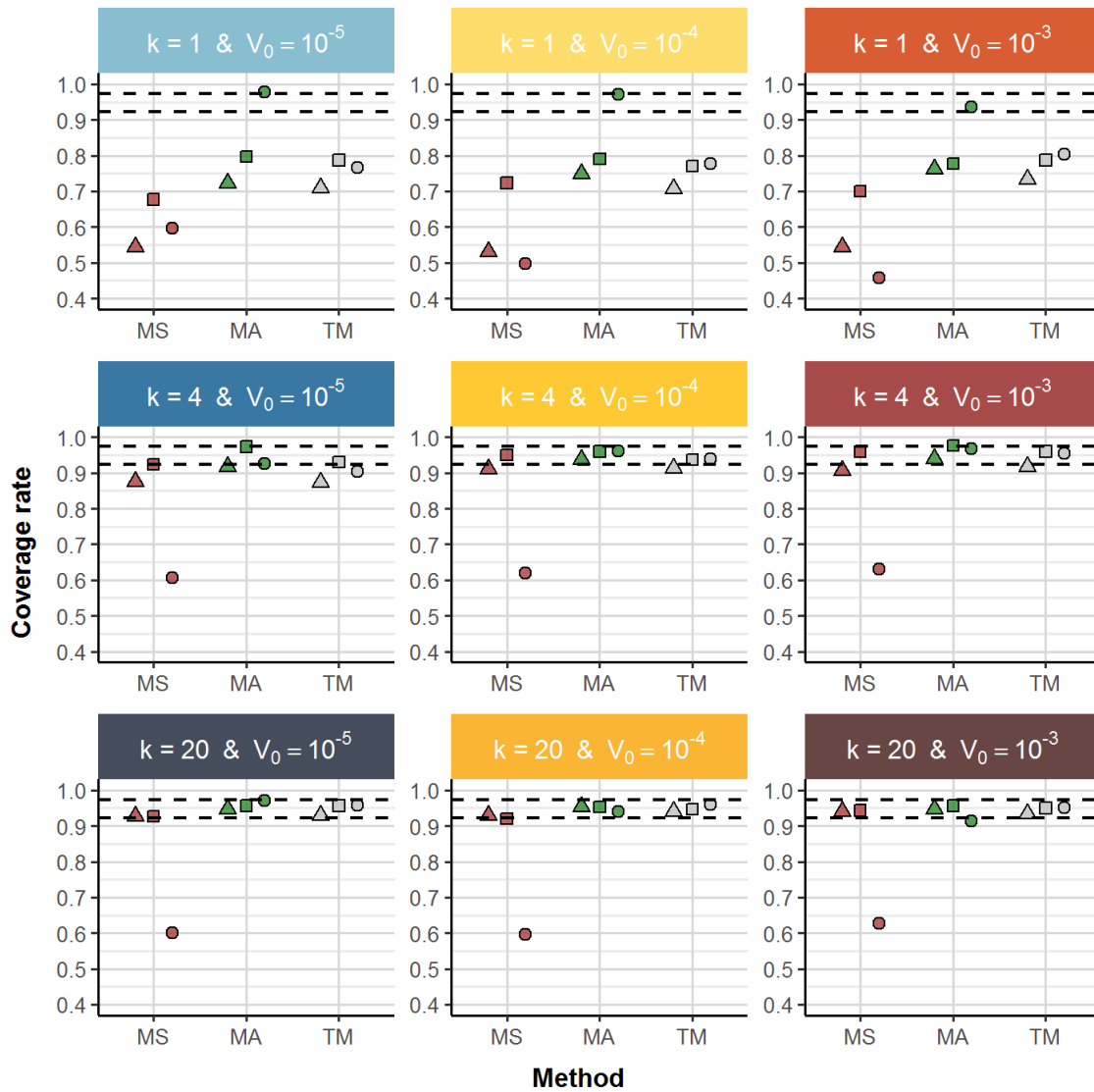
502



503

504 **Figure 2. Setting I.** (A) For each scenario, the percentage of simulations where each
 505 candidate model was selected using AIC. Title of the facet indicates the true model. (B)
 506 Boxplots of weights (whiskers from the 2.5th to the 97.5th percentile) associated with each
 507 candidate model using AIC values. The asterisk denotes the true model in each scenario.

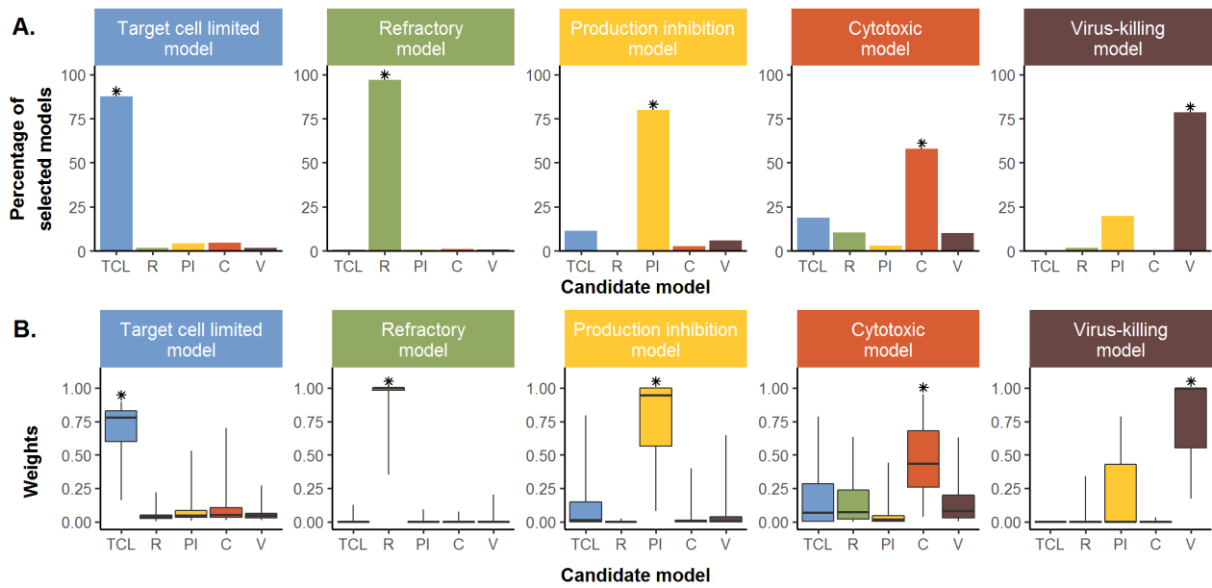
508



509

510 **Figure 3. Coverage rate of R_0 and δ in setting I.** Coverage rate of the parameters R_0 (dots),
 511 π (squares) and δ (triangles) for each scenario using model selection, model averaging or the
 512 true model. Dashed lines represents the prediction interval around 0.95.

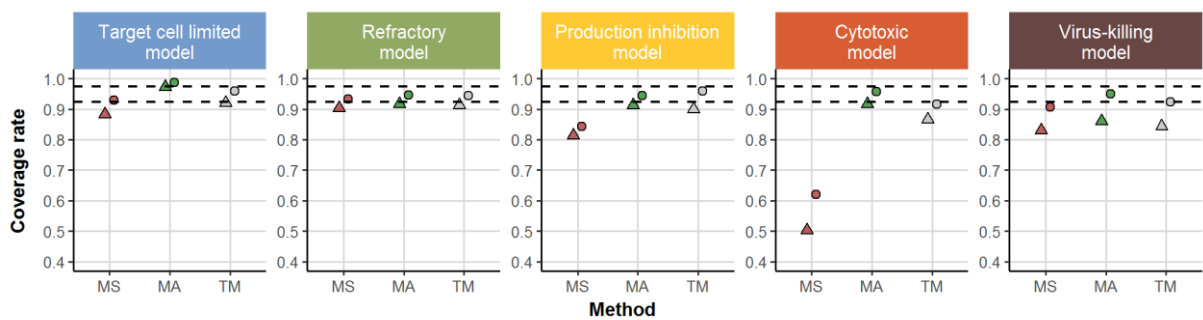
513



514

515 **Figure 4. Setting II.** (A) For each scenario, the percentage of simulations where each
 516 candidate model was selected using AIC. Title of the facet indicates the true model. (B)
 517 Boxplots of weights (whiskers from the 2.5th to the 97.5th percentile) associated with each
 518 candidate model using AIC values. The asterisk denotes the true model in each scenario.

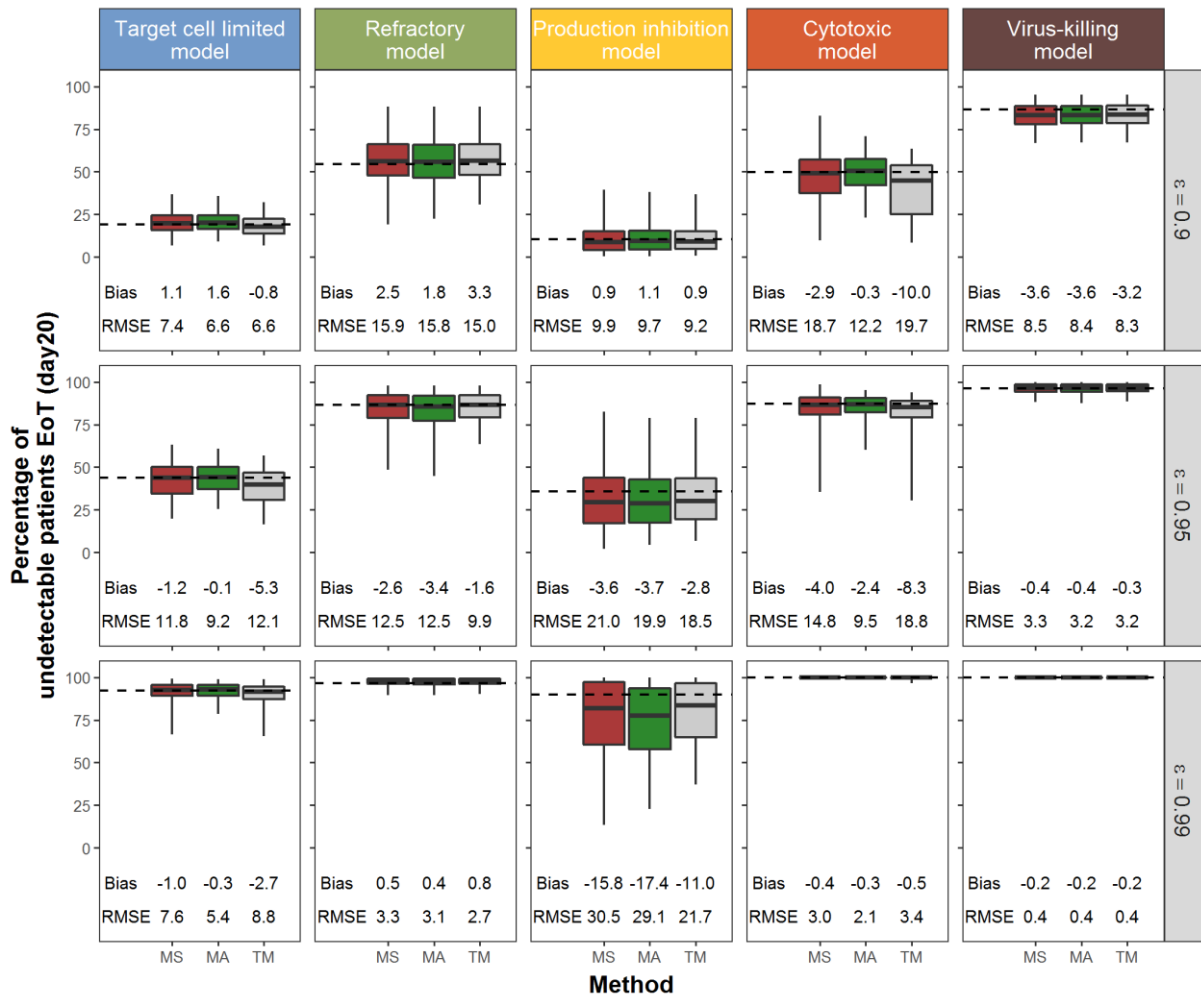
519



520

521 **Figure 5. Coverage rate of R_0 and δ in setting II.** Coverage rate of the parameters R_0 (dots)
 522 and δ (triangles) for each scenario using model selection, model averaging or the true model.
 523 Dashed lines represents the prediction interval around 0.95.

524



525

526 **Figure 6. Distribution of the expected proportion of patients below the limit of detection**
 527 **at day 20 using model selection (red), model averaging (green) or the true model (grey).**

528 For each scenario, whiskers represent the 2.5th to the 97.5th percentile and each row
 529 corresponds to a different value of the treatment effect, noted ϵ .

530

531