

1

2

3

4

5

Calibration of individual-based models to epidemiological data:

6

a systematic review

7

8

9

C. Marijn Hazelbag^{1*}, Jonathan Dushoff^{1,2}, Emanuel M. Dominic¹, Zinhle E. Mthomboti¹,

10

Wim Delva^{1,3,4,5,6,7}

11

12

¹ The South African Department of Science and Technology-National Research Foundation (DST-

13

NRF) South African Centre for Epidemiological Modelling and Analysis (SACEMA), Stellenbosch

14

University, Stellenbosch, South Africa

15

² Department of Mathematics and Statistics, the Institute for Infectious Disease Research,

16

McMaster University, Hamilton, Ontario, Canada

17

³ School for Data Science and Computational Thinking, Stellenbosch University, Stellenbosch,

18

South Africa

19

⁴ Center for Statistics, I-BioStat, Hasselt University, Diepenbeek, Belgium

20

⁵ Department of Global Health, Faculty of Medicine and Health, Stellenbosch University,

21

Stellenbosch, South Africa

22

⁶ International Centre for Reproductive Health, Ghent University, Ghent, Belgium

23

⁷ Rega Institute for Medical Research, KU Leuven, Leuven, Belgium

24

25

* Corresponding author

26

E-mail: marijnhazelbag@sun.ac.za

27 **Abstract**

28 Individual-based models (IBMs) informing public health policy should be calibrated to data and
29 provide estimates of uncertainty. Two main components of model-calibration methods are the
30 parameter-search strategy and the goodness-of-fit (GOF) measure; many options exist for each
31 of these. This review provides an overview of calibration methods used in IBMs modelling
32 infectious disease spread.

33

34 We identified articles on PubMed employing simulation-based methods to calibrate IBMs
35 informing public health policy in HIV, tuberculosis, and malaria epidemiology published
36 between 1 January 2013 and 31 December 2018. Articles were included if models stored
37 individual-specific information, and calibration involved comparing model output to population-
38 level targets. We extracted information on parameter-search strategies, GOF measures, and
39 model validation.

40

41 The PubMed search identified 653 candidate articles, of which 84 met the review criteria. Of the
42 included articles, 40 (48%) combined a quantitative GOF measure with an algorithmic
43 parameter-search strategy – either an optimisation algorithm (14/40) or a sampling algorithm
44 (26/40). These 40 articles varied widely in their choices of parameter-search strategies and GOF
45 measures. For the remaining 44 (52%) articles, the parameter-search strategy could either not
46 be identified (32/44) or was described as an informal, non-reproducible method (12/44). Of
47 these 44 articles, the majority (25/44) were unclear about the GOF measure used; of the rest,
48 only five quantitatively evaluated GOF. Only a minority of the included articles, 14 (17%)
49 provided a rationale for their choice of model-calibration method. Model validation was
50 reported in 31 (37%) articles.

51

52 Reporting on calibration methods is far from optimal in epidemiological modelling studies of
53 HIV, malaria and TB transmission dynamics. The adoption of better documented, algorithmic

54 calibration methods could improve both reproducibility and the quality of inference in model-
55 based epidemiology. There is a need for research comparing the performance of calibration
56 methods to inform decisions about the parameter-search strategies and GOF measures.

57

58

59 **Author summary**

60 Calibration - that is, “fitting” the model to data - is a crucial part of using mathematical models to
61 better forecast and control the population-level spread of infectious diseases. Evidence that the
62 mathematical model is well-calibrated improves confidence that the model provides a realistic
63 picture of the consequences of health policy decisions. To make informed decisions,
64 Policymakers need information about uncertainty: i.e., what is the range of likely outcomes
65 (rather than just a single prediction). Thus, modellers should also strive to provide accurate
66 measurements of uncertainty, both for their model parameters and for their predictions. This
67 systematic review provides an overview of the methods used to calibrate individual-based
68 models (IBMs) of the spread of HIV, malaria, and tuberculosis. We found that less than half of the
69 reviewed articles used reproducible, non-subjective calibration methods. For the remaining
70 articles, the method could either not be identified or was described as an informal, non-
71 reproducible method. Only one-third of the articles obtained estimates of parameter
72 uncertainty. We conclude that the adoption of better-documented, algorithmic calibration
73 methods could improve both reproducibility and the quality of inference in model-based
74 epidemiology.

75

76

77

78

79

80 Introduction

81 Individual-based models (IBMs) intended to inform public health policy should be
82 calibrated to real-world data and provide valid estimates of uncertainty [1, 2]. IBMs track
83 information for a simulated collection of interacting individuals [3]. IBMs allow for more
84 detailed incorporation of heterogeneity, spatial structure, and individual-level adaptation (e.g.
85 physiological or behavioural changes) compared to other modelling frameworks [4]. This
86 complexity makes IBMs valuable planning tools, particularly in settings where real-world
87 intricacies that are not accounted for in simpler models have important effects [5, 6]. However,
88 researchers and policymakers often battle with the question of how much value they can attach
89 to the results of IBMs [7]. Fitting an IBM to empirical data (calibration) improves confidence that
90 the simulation model provides a realistic and accurate estimate of the outcome of health policy
91 decisions (e.g. projection of the disease prevalence under different intervention strategies, or the
92 cost-effectiveness of different intervention strategies) [8-12]. Transparent reporting on
93 calibration methods for IBMs is therefore required [11, 12].

94

95 Parameter values with accompanying confidence intervals used in IBMs are obtained from the
96 literature and are often obtained through statistical estimation. When researchers cannot
97 estimate parameters from empirical data, they obtain their likely values through calibration
98 [12]. Parameter calibration is often difficult for IBMs because their greater complexity can
99 render the likelihood function analytically intractable (i.e. it is impossible to write down the
100 likelihood function in closed form) or prevent explicit numerical calculation of the likelihood
101 function [13-15]. Consequently, simulation-based calibration methods that avoid the use of a
102 likelihood function in closed form have been developed [16]. These methods run the model for
103 different parameter sets to identify parameter sets producing model output that best resembles
104 the summary statistics obtained from the empirical data (e.g. disease prevalence over time).
105 Formal simulation-based calibration requires *summary statistics (targets)* from empirical data, a
106 *parameter-search strategy* for exploring the parameter space, a *goodness-of-fit (GOF)* measure to

107 evaluate the concordance between model output and targets, *acceptance criteria* to determine
108 which parameter sets produce model output close enough to the targets, and a *stopping rule* to
109 determine when the calibrations ends [9,19]. IBMs vary in their complexity (i.e. the number of
110 parameters) and the amount of data available for calibration and validation [10]. Simulation-
111 based calibration of IBMs of higher complexity is typically more computationally intensive [17,
112 18].

113

114 In this review, we pay particular attention to the parameter-search strategy and GOF
115 measure used. Algorithmic parameter-search strategies can be divided into *optimisation*
116 *algorithms* and *sampling algorithms* [14], S2 table describes commonly used algorithms.
117 Optimisation algorithms find the parameter combination that optimises the GOF, resulting in a
118 single best parameter combination. Examples include grid-search and iterative, descent-guided
119 optimisation algorithms using simplex-based or direct search methods (e.g. the Nelder-Mead
120 method) [20], but many different algorithms exist [21]. Optimisation algorithms provide only
121 point estimates of parameters; once these are found, another algorithm may be used to obtain
122 confidence intervals (e.g. the profile likelihood method, Fisher information, etc.) [22, 23].
123 Sampling algorithms aim to find a distribution of parameter values that approximate the
124 likelihood surface or posterior distribution. Examples include approximate Bayesian
125 computation (ABC) methods and sampling importance resampling [8, 13, 14, 24, 25]. Parameter
126 distributions obtained from sampling algorithms allow for the representation of correlations
127 between parameters and for parameter uncertainty to be incorporated into model projections
128 [2, 6, 8, 19, 26]. Quantitative measures of GOF include distance measures (e.g. relative distance,
129 squared distance) and measures based on a surrogate likelihood function: the likelihood of
130 observing the target statistic under the assumption that the model output is a random draw
131 from a presumed distribution (e.g. binomial for prevalence statistics). As the model output is not
132 necessarily distributed as presumed, we refer to this likelihood as the “surrogate” likelihood. A

133 more subjective method of calibration involves the manual adjustment of parameter values,
134 followed by a visual assessment of whether the model outputs resemble empirical data [27].

135

136 Previous research in the context of IBMs of HIV transmission found that 22 (69%) out of
137 32 included articles described the process through which the model was calibrated to data [12].
138 The impact of stochasticity on the model results, defined as the random variation in model
139 output induced by running the model multiple times using the same parameter value with a
140 different random seed, was summarised in nearly half (15/32) of the articles [12]. The depth of
141 reporting on calibration methods was highly variable [9, 12]. A systematic review in the context
142 of population-level health policy models, including 37 articles, found that 25(71%) of these
143 performed model calibration [28]. About half (12/25) of these articles reported on the
144 calibration methods used, whereas the other half (13/25) used informal methods for parameter
145 calibration or did not report on the calibration methods [28]. Previous research on calibration
146 methods in cancer-simulation models in general – not IBMs specifically – found that 131 (85%)
147 out of 154 included articles may have calibrated at least one unknown parameter. Of the 131
148 articles that calibrated parameters, the majority (84/131) did not describe the use of a GOF
149 measure, the rest either used a quantitative GOF (27/131) such as the likelihood or distance
150 measures or used visual assessment of GOF (20/131) [9]. Only a few articles reported parameter
151 distributions resulting from calibration; most only presented a single best parameter
152 combination [9]. Information on the parameter-search strategy and stopping rules was generally
153 not well described, and acceptance criteria were rarely mentioned [9, 29]. Of the 154 articles
154 included in the review by Stout *et al.*, 80 (52%) mentioned model validation [9]. However, while
155 previous studies have reviewed specific portions of the modelling literature, they either did not
156 focus on IBMs or did not focus on the calibration methods in much detail.

157

158 We conducted a systematic review of epidemiological studies using IBMs of the HIV,
159 malaria and tuberculosis (TB) epidemics, as these have been among the most investigated

160 epidemics with the highest global burden of disease [30]. We aim to provide an overview of
161 current practices in the simulation-based calibration of IBMs.

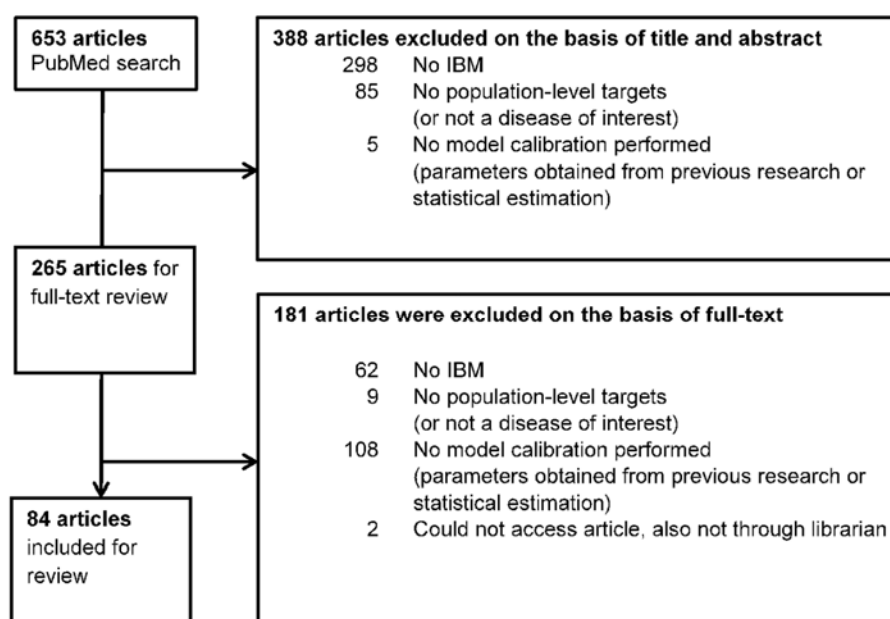
162

163 Results

164 Selection of articles for inclusion

165 The PubMed search resulted in 653 publications, of which 84 articles were included for
166 review; 388 were excluded based on title and abstract, and another 181 were excluded based on
167 a full-text review (see Fig 1). The number of articles selected by publication year increased from
168 seven in 2013 to 20 in 2018.

169



170

171 **Fig 1. PRISMA flow diagram detailing the selection process of articles included in the**
172 **review.**

173 Scope and objectives of included articles

174 S1 Table summarises the characteristics of the included articles. Fifty-eight (69%) of the
175 included articles presented IBMs in HIV research, 16 (19%) concerned malaria, and another 10
176 (12%) concerned tuberculosis.

177

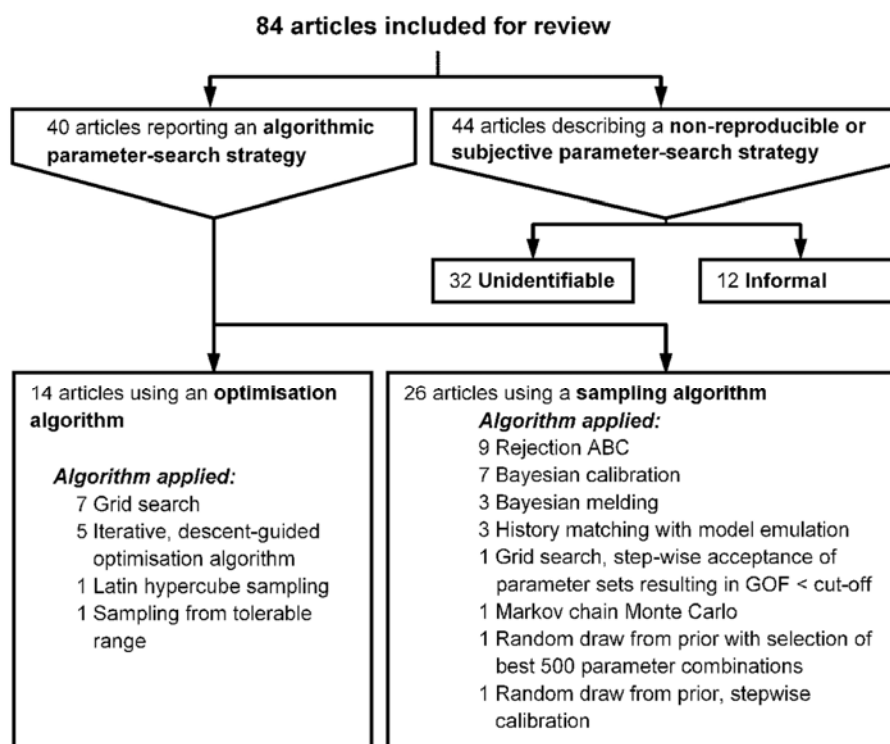
178 Most articles, namely 56 (67%), investigated the effect of an intervention, 17 articles
179 looked at behavioural or biological explanations for the observed epidemic, and other goals (e.g.
180 parameter estimation, model development) were used in 17. In total, six (7%) articles had two
181 objectives. For most of these (5/6), one of the objectives was investigating the effect of an
182 intervention (see S1 Table).

183

184 **Parameter-search strategies and measures of GOF**

185 Of the included articles, 40 (48%) combined a quantitative measure of GOF with an
186 algorithmic parameter-search strategy, which was an optimisation algorithm (14/40) or a
187 sampling algorithm (26/40) (see Fig 2). For the remaining 44 (52%) articles, the parameter-
188 search strategy could either not be identified (32/44) or was described as an informal, non-
189 reproducible method (12/44). Tables A, B and C in S1 appendix show that there is no convincing
190 evidence that the parameter search strategy changed with publication year or differed by
191 disease studied. A brief description of the methods referred to in Fig 2 under optimisation
192 algorithm and sampling algorithm is provided in S2 Table.

193



194

195 **Fig 2. Reporting and application of parameter search strategies in epidemiological**
196 **studies.**

197

198 Detailed information on calibration methods for the 14 (17%) articles using optimisation
199 algorithms is reported in Table 1. For the parameter-search strategy, most articles used either a
200 grid search (7/14), Latin square (1/14) or random draw from tolerable range (1/14), followed
201 by the selection of the single best parameter combination. Several iterative, descent-guided
202 optimisation algorithms (i.e. Nelder-Mead, interior-point algorithm, coordinate descent with
203 golden section search, random search mechanism) were used in the remaining articles (5/14).
204 Of these five articles, most (4/5) accepted a single best parameter combination without
205 confidence intervals, while the remaining article obtained confidence intervals around
206 parameter estimates (see S1 Text.). For the GOF measure, the most common choice was a
207 squared distance (6/14). Various GOF measures were used in the remaining articles; these
208 include absolute distances (2/14) and R-squared (2/14).

209

210

211 **Table 1. Details of the calibration methods used in articles using optimisation algorithms**
 212 **for calibration, sorted by parameter search strategy algorithm.**

Authors	Year	Pathogen	Parameter search strategy algorithm	GOF
Luo <i>et al.</i>	2018	HIV	Grid search	Absolute distance
Romero-Severson <i>et al.</i>	2013	HIV	Grid search	Kolmogorov-Smirnov
Marshall <i>et al.</i>	2018	HIV	Grid search	R-squared
Goedel <i>et al.</i>	2018	HIV	Grid search	R-squared and Manhattan distance of parameters
Brookmeyer <i>et al.</i>	2014	HIV	Grid search	Squared distance
Suen <i>et al.</i>	2014	TB	Grid search	Number of model outputs within the confidence intervals around the targets
Suen <i>et al.</i>	2015	TB	Grid search	Number of model outputs within the confidence intervals around the targets
Bershteyn <i>et al.</i>	2013	HIV	Iterative, descent-guided optimisation algorithm (<i>Coordinate descent w. golden section search</i>)	Squared distance
Klein <i>et al.</i>	2015	HIV	Iterative, descent-guided optimisation algorithm (<i>Coordinate descent w. golden section search</i>)	Squared distance
Sauboin <i>et al.</i>	2015	Malaria	Iterative, descent-guided optimisation algorithm (<i>Interior point algorithm, hill-climbing</i>)	Squared distance
Knight <i>et al.</i>	2015	TB, HIV	Iterative, descent-guided optimisation algorithm (<i>Nelder-Mead</i>)	Squared distance
Kasaie <i>et al.</i>	2018	HIV	Iterative, descent-guided optimisation algorithm (<i>Random search mechanism</i>)	Absolute distance
Shrestha <i>et al.</i>	2017	TB	Latin hypercube sampling	Surrogate likelihood
Jewell <i>et al.</i>	2015	HIV	Sampling from tolerable range	Squared distance

213

214 Table 2 contains the details of the calibration methods in the 26 (31%) articles using
 215 sampling algorithms. Random sampling from the prior, followed by rejection ABC, was used the
 216 most (9/26). Bayesian calibration (7/26), Bayesian melding (3/26) and history matching with
 217 model emulation (3/26) were also used. Most articles (10/26) used the surrogate likelihood as a
 218 measure of GOF, and Various GOF measures were used in the remaining articles, these include
 219 absolute distances (4/26), relative distances (4/26) and squared distances (4/26). (see Table 2).

220

221 **Table 2. Details of the calibration methods in articles using sampling algorithms for**
 222 **calibration, sorted by parameter search strategy algorithm.**

Authors	Year	Pathogen	Parameter search strategy algorithm	GOF
Cameron <i>et al.</i>	2015	Malaria	Bayesian calibration (<i>Combining model emulation with MCMC</i>)	Surrogate likelihood
Huynh <i>et al.</i>	2015	TB	Bayesian calibration (<i>Latin hypercube with IMIS</i>)	Surrogate likelihood
Chang <i>et al.</i>	2018	TB	Bayesian calibration (<i>Latin hypercube with IMIS</i>)	Surrogate likelihood
Penny <i>et al.</i>	2015	Malaria	Bayesian calibration (<i>MCMC</i>)	Surrogate likelihood
Penny <i>et al.</i>	2015	Malaria	Bayesian calibration (<i>MCMC</i>)	Surrogate likelihood
White <i>et al.</i>	2018	Malaria	Bayesian calibration (<i>MCMC</i>)	Surrogate likelihood
Schalkwyk <i>et al.</i>	2018	HIV	Bayesian calibration (<i>Random draw from prior with SIR</i>)	Surrogate likelihood
Abuelezam <i>et al.</i>	2016	HIV	Bayesian melding	Squared distance
McCormick <i>et al.</i>	2014	HIV	Bayesian melding	Surrogate likelihood
McCormick <i>et al.</i>	2017	HIV	Bayesian melding	Surrogate likelihood
Ciaranello <i>et al.</i>	2013	HIV	Grid search, step-wise acceptance of parameter sets resulting in GOF < cut-off	Absolute distance
McCreesh <i>et al.</i>	2017	HIV	History matching with model emulation	Implausibility measure
McCreesh <i>et al.</i>	2017	HIV	History matching with model emulation	Implausibility measure
McCreesh <i>et al.</i>	2018	HIV	History matching with model emulation	Implausibility measure
Shcherbacheva <i>et al.</i>	2018	Malaria	Markov chain Monte Carlo	Absolute distance
Johnson <i>et al.</i>	2016	HIV	Random draw from prior with selection of best 500 parameter combinations	Surrogate likelihood
Pizzitutti <i>et al.</i>	2015	Malaria	Random draw from prior, stepwise calibration	Absolute distance
Nakagawa <i>et al.</i>	2016	HIV	Rejection ABC (<i>Random draw from prior</i>)	Relative distance
Nakagawa <i>et al.</i>	2017	HIV	Rejection ABC (<i>Random draw from prior</i>)	Chi-square
Pizzitutti <i>et al.</i>	2018	Malaria	Rejection ABC (<i>Random draw from prior</i>)	Squared distance
Cambiano <i>et al.</i>	2018	HIV	Rejection ABC (<i>Random draw from prior</i>)	Relative distance
Hontelez <i>et al.</i>	2013	HIV	Rejection ABC (<i>Random draw from prior</i>)	Squared distance
Phillips <i>et al.</i>	2013	HIV	Rejection ABC (<i>Random draw from prior</i>)	Relative distance
Phillips <i>et al.</i>	2015	HIV	Rejection ABC (<i>Random draw from prior</i>)	Relative distance
Shrestha <i>et al.</i>	2017	HIV	Rejection ABC (<i>Random draw from prior</i>)	Absolute distance
Tuite <i>et al.</i>	2017	TB	Rejection ABC (<i>Random draw from prior</i>)	Squared distance

223 IMIS, Incremental-mixture importance sampling; SIR, Sampling importance resampling; MCMC, Markov chain Monte
 224 Carlo.

225

226 From the 44 (52%) articles with unidentifiable or informal parameter-search strategies,
227 the majority (25/44) are also unclear about the GOF used, while the rest either relied on visual
228 inspection as a GOF (14/44) or used a quantitative GOF (5/44).

229

230 Only 14 (17%) of the 84 included articles provided a rationale for their choice of model-
231 calibration method. For example, McCreesh *et al.* [32] reported: “The model was fitted to the
232 empirical data using history matching with model emulation, which allowed uncertainties in
233 model inputs and outputs to be fully represented, and allowed realistic estimates of uncertainty
234 in model results to be obtained” (see S2 Text. for more examples). Other examples indicate that
235 an algorithmic calibration method failed to provide either a good fit or parameter estimates:
236 “Ultimately, we chose to use visual inspection because the survival curves did not fit closely
237 enough using the other two more quantitative approaches.” [33] Or “[Calibration] was unable to
238 resolve co-varying parameters. These parameters were adjusted by hand...” [34].

239 Ten out of the 84 articles included (12%) used a weighted calculation of GOF. Four
240 articles weighted the GOF based on the amount of data behind the summary statistic fitted to, for
241 example by weighting based on the inverse of the width of the confidence interval around the
242 data. In contrast, one article increased the weight for a data source for which fewer data was
243 available. Other strategies included weighting based on a subjective assessment of the quality of
244 the data, or weighting based on which data they wanted the model to fit best. One article down-
245 weighted particular data to improve fit. Others stressed the importance of determining weights
246 a priori since weights are chosen subjectively.

247

248 **Acceptance criteria and stopping rules**

249 None (0/14) of the articles applying optimisation algorithms mentioned the acceptance
250 criteria or stopping rules. Acceptance criteria and stopping rules applied in studies using
251 sampling algorithms can be summarised as running the model until obtaining an arbitrary
252 number of accepted parameter combinations.

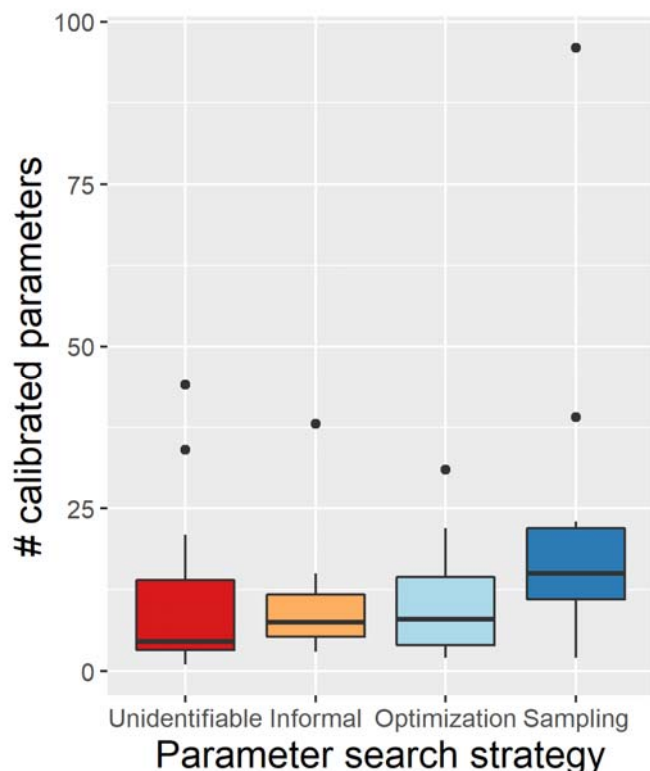
253

254 **The number of target statistics, the number of calibrated parameters and the size**
255 **of the simulated population**

256 The number of target statistics was explicitly mentioned in only three (4%) of the 84
257 included articles, for 62 (74%) articles we had enough information to attempt to deduce this
258 number from either text or figures. The remaining 19 (23%) articles either provided incomplete
259 information (11/19) or no information (8/19). Some (4/65) of the articles for which we were
260 able to obtain the number of target statistics had different numbers of target statistics for
261 calibration in different locations or calibration to different diseases. The 61 (73%) articles for
262 which we were able to obtain a single count had a median number of target statistics of 23
263 (range 1 - 321). A histogram of the number of target statistics is provided in figure A in S2
264 Appendix.

265 The number of calibrated parameters was explicitly mentioned in 11 (13%) of the 84
266 included articles, for another 53 (63%) articles it was possible to deduce this number from
267 either text or figures. For the remaining 20 (24%) articles either provided incomplete
268 information (10/20) or no information at all (10/20). The 64 (75%) articles for which we were
269 able to obtain a count had a median number of calibrated parameters of 10 (range 1 - 96). A
270 histogram of the number of calibrated parameters is provided in figure B in S2 Appendix. A
271 Kruskal-Wallis test was performed to examine differences in the number of calibrated
272 parameters among parameter search strategies (See Fig 3). Significant differences (Chi-square =
273 9.304, $p = 0.026$, $df = 3$) were found amid the four categories of parameter search strategies.
274 Pairwise comparisons using the Wilcoxon rank-sum test corrected for multiple testing by
275 Benjamini-Hochberg correction showed that articles using Sampling strategies have higher
276 numbers of calibrated parameters compared to papers for which we could not identify the
277 parameter search strategy.

278



279

280 **Fig 3. Boxplots of the number of target statistics by the parameter search strategy.**

281

282 For 55 (66%) articles, we obtained counts for both the number of target statistics and
283 the number of calibrated parameters. For many of these articles (17/55), the number of
284 calibrated parameters appeared to exceed the number of target statistics. A plot of the number
285 of target statistics against the number of calibrated parameters is provided in figure C in S2
286 Appendix.

287 The size of the simulated population was explicitly mentioned in 54 (64%) of the 84
288 included articles, for another 9 (11%) articles it was possible to deduce this number from either
289 text or figures. The remaining 21 (25%) articles either provided incomplete information (3/21)
290 or no information at all (18/21). For the 63 (75%) articles for which we obtained a number, the
291 median population size was 78000 (range: 250 - 47000000). A histogram of the \log_{10} of the size
292 of the simulated population is provided in figure D in S2 Appendix.

293

294

295 **Computational aspects and the use of platforms**

296 The software used to build IBM was not reported in 33 (39%) of the articles. Sixteen
297 articles (19%) used the low-level programming language C++, six (7%) used MATLAB, and
298 another six (7%) used Python. Various other computing platforms were used in the remaining
299 23 (28%) articles. A high-performance computing facility was used in 16 (19%) articles.

300

301 Several simulation tools (i.e. CEPAC [35], EMOD [36] HIV-CDM [37], MicroCOSM [38],
302 PATH [39], STDSIM [40] and TITAN [41]) were used in the articles modelling HIV. Similarly, two
303 platforms (i.e. EMOD [42] and OpenMalaria [43]) were used in the articles modelling malaria. In
304 the articles modelling tuberculosis, the only tool reported was EMOD [44].

305

306 **Model validation**

307 Only 31 (37%) articles mentioned that a validation of the model had been performed.

308

309 **Discussion**

310 More than half of IBMs we studied used non-reproducible or subjective calibration methods.
311 Articles that reported the use of formal calibration methods used a wide range of parameter-
312 search strategies and GOF measures. Only one-third of articles used calibration methods that
313 quantify parameter uncertainty. These findings are important because choices concerning the
314 calibration method can have substantial effects on model results and policy implications [2, 6-8,
315 45-47].

316

317 We encourage authors to use the standardised Calibration Reporting Checklist of Stout *et*
318 *al.* [9]. While algorithmic parameter-search strategies are in principle reproducible, unclear or
319 incomplete reporting, and non-disclosure of software code can render them de facto non-

320 reproducible. [48]. Manual adjustment of parameter values and visual inspection of GOF may
321 perform equally well compared to other methods in terms of GOF alone [49], may provide
322 researchers with valuable insights into and familiarity with the model [50], and can be useful for
323 purely didactic purposes [51-53]. However, we advise against using these methods in analyses
324 intended to inform public health as they do not favour reproducibility and involve subjective
325 judgment, which may produce less than optimal calibration results and usually leads to the
326 acceptance of a single parameter set (i.e. does not provide parameter uncertainty) [19]. On
327 occasion, authors justified their choice of an informal method by indicating that algorithmic
328 calibration methods did not converge to provide parameter estimates, or failed to provide a
329 satisfactory fit to the targets. A potential explanation for non-convergence of an algorithmic
330 calibration method is that the parameters in question are unidentifiable, which is the case when
331 a vast array of different parameter combinations provide a comparably good fit to the target
332 statistics. Performing manual calibration in such an instance will deliver one set of parameters
333 out of all of the parameter combinations that provide a fit. However, using this single parameter
334 combination hides the fact that there is not enough information to uniquely identify the best
335 parameter values. Furthermore, model-stochasticity provides the possibility that a great fit is
336 found by chance for a parameter combination for which the probability of observing the target
337 statistics is lower than for other parameter combinations.

338

339 There are several methodological challenges in the calibration of individual-based
340 models, including the choice of calibration method – i.e. the combination of algorithmic
341 parameter-search strategy and GOF measure. The findings of the current review and previous
342 research suggest that there is no consensus on which calibration method to use [9, 10, 19, 54,
343 55]. Additionally, some of the articles reviewed here indicated that algorithmic calibration
344 methods had failed, leading the researchers to calibrate the model, either fully or partially, by
345 hand. These issues suggest that there is a need for research comparing the performance of
346 calibration methods to inform decisions about the choice of parameter-search strategy and GOF
347 [10]. Previous research on calibration methods focused on the GOF [27], computation time and

348 analyst time [49]. Where applicable, correct estimation of the posterior [56] should be a core
349 aspect of performance. We further suggest investigating several contextual variables, including
350 the amount and nature of the empirical data to calibrate against, the number and type of model
351 parameters to be calibrated and insights to be derived from the calibrated model. As evident
352 from our review, these contextual variables vary widely across IBM studies in epidemiology.

353

354 Another methodological challenge in the calibration of IBMs is determining a priori
355 whether the target statistics provide sufficient information to calibrate the parameters [57],
356 especially when the model has many parameters [58]. Firstly, the target statistics are based on
357 variable amounts of raw data. Secondly, a time series of target statistics is often used, typically
358 violating the assumption of independence implied by many calibration methods. Thirdly, the
359 complexity of the model may hamper an appropriate specification of a prior parameter-
360 distribution (including the specification of a correlation between parameters) that is fully
361 informed by prior knowledge of the data-generating processes represented by the model. These
362 problems preclude the use of standard statistical methods for calculating the number of target
363 statistics that is sufficient for parameter calibration. A related problem is that target summary
364 statistics are based on data from different sources, including observational data that are
365 potentially affected by treatment-confounder feedback (e.g. time-dependent confounder CD4
366 cell count affected by prior cART treatment) [59]. Another related problem is that of validation,
367 i.e. testing model performance on data that was not included in the calibration step. There is
368 considerable debate on when data should be reserved for this purpose [55].

369

370 The last methodological aspect of IBMs we would like to draw attention to is the size of
371 the simulated population [1, 60]. Intuitively, one would recommend that the simulated
372 population size should be similar to the size of the population from which the samples were
373 drawn that gave rise to the target statistics. However, for many studies, modelling the full
374 population is not feasible with currently available computational infrastructure. Instead,
375 researchers often adjust for the inflated stochasticity in the modelled system by averaging

376 outcomes of interest over multiple simulations runs per parameter set [60]. How choices around
377 modelled population size and analysis of model output affect the validity of model inference
378 deserves further attention in future research.

379

380 Our results in the setting of HIV, TB and malaria IBMs indicate that the use of formal
381 calibration methods (48% of articles) is higher than in previous research on simulation models
382 in general – not IBMs specifically. Previously, only one-fifth to one-third of articles reporting on
383 epidemiological models used a quantitative GOF [9, 61]. Our results concerning parameter
384 uncertainty are also optimistic compared to previous research by Stout *et al.* on calibration
385 methods in cancer models, which found that almost no articles quantified parameter
386 uncertainty, but instead accepted a single best-fitting parameter set as the result of the
387 calibration [9]. The same researchers reported that several different combinations of parameter-
388 search strategies and GOFs were used [9], outcomes which are similar to our findings. Stout *et al.*
389 report that articles rarely describe acceptance criteria and stopping rules. Stout *et al.* also report
390 that a standard description of the calibration process lacks in almost all articles [9]. Similarly,
391 previous research on IBMs of HIV transmission found that reporting was lacking in the
392 description of calibration methods [12]. All of this is in agreement with the results of the current
393 review. Concerning the goals of the included articles, our results broadly agree with
394 Punyacharoensin *et al.* They found that the main goals of HIV transmission models for the study
395 of men who have sex with men are: making projections for the epidemic, investigating how the
396 incorporation of various assumptions around the behavioural or biological characteristics affect
397 these projections, and evaluating the impact of interventions [61].

398

399 To our knowledge, this is the first detailed review of methods used to calibrate IBMs of
400 HIV, malaria and TB epidemics. A limitation of our study is that we are unsure to what extent the
401 results are generalisable to other infectious diseases. We encourage future research on other
402 diseases to confirm or refute our current findings on the use of and reporting on methods in the
403 calibration of IBMs in epidemiological research. Similarly, since our PubMed search excluded

404 articles matching “molecular”, we may have missed relevant articles. However, we don’t believe
405 this selection is likely to bias the findings of this review. Another possible concern is that we
406 don’t control for overlaps in authorship; thus, we effectively treat articles that come from a given
407 “research group” as independent observations, even though the calibration method used by a
408 particular group is often the same, as we show in Tables 1 and 2. Another limitation is that the
409 counts presented in this review often had to be deduced from the article, this was a difficult and
410 laborious task involving manual counting of target statistics in either the text, figures or tables, a
411 process that is prone to error. A final limitation is that we did not go into the strengths and
412 weaknesses of each method. Existing literature compares the performance of alternative
413 algorithms for calibrating the same model but does not allow us to draw general conclusions
414 [10]. As a starting point for comparison, we provide a brief description of calibration methods in
415 S2 Table.

416

417 In conclusion, it appears that calibrating individual-based models in epidemiological studies of
418 HIV, malaria and TB transmission dynamics remains more of an art than a science. Besides
419 limited reproducibility for a majority of the modelling studies in our review, our findings raise
420 concerns over the correctness of model inference (e.g., estimated impact of past or future
421 interventions) for models that are poorly calibrated. The quality of inference and reproducibility
422 in model-based epidemiology could benefit from the adoption of algorithmic parameter-search
423 strategies and better-documented calibration and validation methods. We recommend the use of
424 sampling algorithms to obtain valid estimates of parameter uncertainty and correlations
425 between parameters. There is a need for simulation-based studies that compare the
426 performance, strengths and limitations of different methods for calibrating IBMs to
427 epidemiological data.

428

429 **Materials and methods**

430 This review was performed following the Preferred Reporting Items for Systematic
431 Reviews and Meta-Analyses (PRISMA) statement [31]. The PRISMA flow diagram details the
432 selection process of articles included for review (see Fig 1).

433

434 **Search strategy and selection criteria**

435 We identified articles on PubMed that employed simulation-based methods to calibrate
436 IBMs of HIV, malaria and tuberculosis, and that were published between 1 January 2013 and 31
437 December 2018. Six years seemed to be long enough to yield a sizeable amount of information
438 and to potentially observe any (recent) time trends, and short enough to be feasible and to speak
439 to recent practices in model calibration in epidemiological modelling studies. The following
440 search query was performed on 31 January 2019: *'((HIV[tiab] OR malaria[tiab] OR*
441 *tuberculo*[tiab] OR TB[tiab]) AND (infect* OR transmi* OR prevent*) AND (computer*
442 *simulation[tiab] OR microsimulation[tiab] OR simulation[tiab] OR agent-based[tiab] OR*
443 *individual-based[tiab] OR computer model*[tiab] OR computerized model*[tiab]) AND*
444 *("2013/01/01"[Date - publication] : "2018/12/31"[Date - publication]) NOT(molecular))'*.

445

446 Eligibility criteria were agreed upon by WD, JD and CMH before screening. Articles were
447 included if models stored individual-specific information and calibration involved running the
448 model and comparing model output to population-level targets expressed as summary statistics.
449 We excluded review articles, statistical simulation studies, and studies that focused on molecular
450 biology and immunology because we were primarily interested in studies informing public
451 health policy.

452

453 Titles and abstracts were screened for eligibility by CMH, and difficult cases were
454 discussed with WD. If the title and abstract did not provide sufficient information for exclusion, a

455 full-text examination was performed. Full-text inclusion was performed by two independent
456 researchers (CMH and either ZM or ED) for a subset of 100 articles. CMH included 28 articles, of
457 which ZM and ED did not include six; these six articles were double-checked by WD and
458 consequently included for review. ZM included four articles that CMH did not include; these four
459 articles were double-checked by WD and consequently not included for review. After that, full-
460 text inclusion was performed by CMH in consultation with WD.

461

462 **Data extraction**

463 For each article, we extracted information on the objective of the study (i.e. estimating
464 the effect of an intervention, investigating a behavioural or biological explanation for the
465 observed infectious disease outbreak or other goals including estimation of parameters or
466 model development), the parameter-search strategy and the GOF measure, the rationale for
467 choosing this calibration strategy over alternatives, and model validation. Acceptance criteria
468 and stopping rules are only relevant for articles applying algorithmic parameter-search
469 strategies and collected for that subset of articles. For readability purposes, we say “used” to
470 mean “reported the use of” throughout this review.

471

472 Information was collected independently by two reviewers (CMH and either ZM or ED)
473 for each article included using a prospectively developed form. This form was based on the
474 Calibration Reporting Checklist of Stout *et al.* [9] and was extended by several items, including;
475 the software and hardware used to build the model, the size of the initial population of agents
476 and the name of the modelling platform. Additionally, we inserted several items to collect
477 information on the number of calibrated parameters, the number of fixed parameters, and the
478 number of targets. We noted how information on these counts was reported in the articles (i.e.
479 the number was explicitly provided, could be deduced from text or figures, was provided
480 incompletely or was not provided).

481

482 Information on calibration methods was extracted verbatim, allowing for later
483 classification. Articles on which there was disagreement in the classification were discussed by
484 WD, JD and CMH until an agreement was reached. We classified articles reporting both
485 algorithmic and informal calibration as informal since doing part of the calibration informally
486 makes the entire calibration irreproducible.

487

488 **Acknowledgements**

489 The authors gratefully acknowledge the help of all SACEMA students and researchers,
490 specifically the fruitful conversations and helpful comments on the manuscript by Prof. Alex
491 Welte, Mrs Cari van Schalkwyk, Dr Florian Marx, Prof. Juliet Pulliam and Dr Larisse Bolton. We
492 would also like to acknowledge Mrs Marisa Honey and Mrs Susan Lotz from the Stellenbosch
493 writing lab, who copy-edited a first version of the manuscript.

494

495 **References**

- 496 [1] Bobashev GV, Morris RJ. Uncertainty and inference in agent-based models. In: 2010
497 Second International Conference on Advances in System Simulation. IEEE; 2010. p. 67–71.
- 498 [2] Briggs AH, Weinstein MC, Fenwick EA, Karnon J, Sculpher MJ, Paltiel AD. Model
499 parameter estimation and uncertainty analysis: a report of the ISPOR-SMDM Modeling Good
500 Research Practices Task Force Working Group–6. Medical decision making. 2012;32(5):722–
501 732.
- 502 [3] Willem L, Verelst F, Bilcke J, Hens N, Beutels P. Lessons from a decade of individual-
503 based models for infectious disease transmission: a systematic review (2006-2015). BMC
504 infectious diseases. 2017;17(1):612.
- 505 [4] Hammond RA. Considerations and best practices in agent-based modeling to inform
506 policy. In: Assessing the use of agent-based models for tobacco regulation. National Academies
507 Press (US); 2015. .

- 508 [5] Johnson LF, Geffen N. A comparison of two mathematical modeling frameworks for
509 evaluating sexually transmitted infection epidemiology. *Sexually transmitted diseases*.
510 2016;43(3):139–146.
- 511 [6] Kennedy MC, O’Hagan A. Bayesian calibration of computer models. *Journal of the Royal
512 Statistical Society: Series B (Statistical Methodology)*. 2001;63(3):425–464.
- 513 [7] Egger M, Johnson L, Althaus C, Schöni A, Salanti G, Low N, et al. Developing WHO
514 guidelines: Time to formally include evidence from mathematical modelling studies.
515 *F1000Research*. 2017;6:1584.
- 516 [8] Menzies NA, Soeteman DI, Pandya A, Kim JJ. Bayesian methods for calibrating health
517 policy models: a tutorial. *Pharmacoeconomics*. 2017;35(6):613–624.
- 518 [9] Stout NK, Knudsen AB, Kong CY, McMahon PM, Gazelle GS. Calibration methods used in
519 cancer simulation models and suggested reporting guidelines. *Pharmacoeconomics*.
520 2009;27(7):533–545.
- 521 [10] Dahabreh IJ, Chan JA, Earley A, Moorthy D, Avendano EE, Trikalinos TA, et al. A Review of
522 Validation and Calibration Methods for Health Care Modeling and Simulation. In: *Modeling and
523 Simulation in the Context of Health Technology Assessment: Review of Existing Guidance, Future
524 Research Needs, and Validity Assessment [Internet]*. Agency for Healthcare Research and Quality
525 (US); 2017. .
- 526 [11] Caro JJ, Eddy DM, Kan H, Kaltz C, Patel B, Eldessouki R, et al. Questionnaire to assess
527 relevance and credibility of modeling studies for informing health care decision making: an
528 ISPOR-AMCP-NPC Good Practice Task Force report. *Value in health*. 2014;17(2):174–182.
- 529 [12] Abuelezam NN, Rough K, Seage III GR. Individual-based simulation models of HIV
530 transmission: reporting quality and recommendations. *PloS one*. 2013;8(9):e75624.
- 531 [13] Lintusaari J, Gutmann MU, Dutta R, Kaski S, Corander J. Fundamentals and recent
532 developments in approximate Bayesian computation. *Systematic biology*. 2017;66(1):e66–e82.

- 533 [14] Hartig F, Calabrese JM, Reineking B, Wiegand T, Huth A. Statistical inference for
534 stochastic simulation models–theory and application. *Ecology letters*. 2011;14(8):816–827.
- 535 [15] Busetto AG, Buhmann JM. Stable Bayesian parameter estimation for biological dynamical
536 systems. In: 2009 International Conference on Computational Science and Engineering. vol. 1.
537 IEEE; 2009. p. 148–157.
- 538 [16] Leombruni R, Richiardi M. Why are economists sceptical about agent-based simulations?
539 *Physica A: Statistical Mechanics and its Applications*. 2005;355(1):103–109.
- 540 [17] Sun NZ, Sun A. Model calibration and parameter estimation: for environmental and
541 water resource systems. Springer; 2015.
- 542 [18] Bellman R. Dynamic programming. Princeton, USA: Princeton University Press.
543 1957;1(2):3.
- 544 [19] Vanni T, Karnon J, Madan J, White RG, Edmunds WJ, Foss AM, et al. Calibrating models in
545 economic evaluation. *Pharmacoeconomics*. 2011;29(1):35–49.
- 546 [20] Nelder JA, Mead R. A simplex method for function minimization. *The computer journal*.
547 1965;7(4):308–313.
- 548 [21] Amaran S, Sahinidis NV, Sharda B, Bury SJ. Simulation optimization: a review of
549 algorithms and applications. *Annals of Operations Research*. 2016;240(1):351–380.
- 550 [22] Joshi M, Seidel-Morgenstern A, Kremling A. Exploiting the bootstrap method for
551 quantifying parameter confidence intervals in dynamical systems. *Metabolic engineering*.
552 2006;8(5):447–455.
- 553 [23] Stryhn H, Christensen J. Confidence intervals by the profile likelihood method, with
554 applications in veterinary epidemiology. In: Proceedings of the 10th International Symposium
555 on Veterinary Epidemiology and Economics, Vina del Mar; 2003. p. 208.

- 556 [24] McKinley T], Vernon I, Andrianakis I, McCreesh N, Oakley JE, Nsubuga RN, et al.
557 Approximate Bayesian Computation and simulation-based inference for complex stochastic
558 epidemic models. *Statistical science*. 2018;33(1):4–18.
- 559 [25] D R. Using the SIR algorithm to simulate posterior distributions. *Bayesian Stat*.
560 1988;3:395–402.
- 561 [26] Poole D, Raftery AE. Inference for deterministic simulation models: the Bayesian melding
562 approach. *Journal of the American Statistical Association*. 2000;95(452):1244–1255.
- 563 [27] Schunn CD, Wallach D, et al. Evaluating goodness-of-fit in comparison of models to data.
564 *Psychologie der Kognition: Reden and vorträge anlässlich der emeritierung von Werner Tack*.
565 2005;p. 115–154.
- 566 [28] Conrads-Frank A, Jahn B, Bundo M, Sroczynski G, Mühlberger N, Bicher M, et al. A
567 Systematic Review Of Calibration In Population Models. *Value in Health*. 2017;20(9):A745.
- 568 [29] Afzali HHA, Gray J, Karnon J. Model performance evaluation (validation and calibration)
569 in model-based studies of therapeutic interventions for cardiovascular diseases. *Applied health
570 economics and health policy*. 2013;11(2):85–93.
- 571 [30] Furuse Y. Analysis of research intensity on infectious disease by disease burden reveals
572 which infectious diseases are neglected by researchers. *Proceedings of the National Academy of
573 Sciences*. 2019;116(2):478–483.
- 574 [31] Moher D, Liberati A, Tetzlaff J, Altman DG. Preferred reporting items for systematic
575 reviews and meta-analyses: the PRISMA statement. *Annals of internal medicine*.
576 2009;151(4):264–269.
- 577 [32] McCreesh N, Andrianakis I, Nsubuga RN, Strong M, Vernon I, McKinley T], et al. Universal
578 test, treat, and keep: improving ART retention is key in cost-effective HIV control in Uganda.
579 *BMC infectious diseases*. 2017;17(1):322.

- 580 [33] Kessler J, Nucifora K, Li L, Uhler L, Braithwaite S. Impact and Cost-Effectiveness of
581 Hypothetical Strategies to Enhance Retention in Care within HIV Treatment Programs in East
582 Africa. *Value in health : the journal of the International Society for Pharmacoeconomics and*
583 *Outcomes Research*. 2015 dec;18(8):946–955. Available from: [http://linkinghub.elsevier.com/-](http://linkinghub.elsevier.com/retrieve/pii/S1098301515050731)
584 [retrieve/pii/S1098301515050731](http://linkinghub.elsevier.com/retrieve/pii/S1098301515050731).
- 585 [34] Klein DJ, Eckhoff PA, Bershteyn A. Targeting HIV services to male migrant workers in
586 southern Africa would not reverse generalized HIV epidemics in their home communities: A
587 mathematical modeling analysis. *International Health*. 2015 mar;7(2):107–113.
- 588 [35] Walensky RP, Borre ED, Bekker LG, Resch SC, Hyle EP, Wood R, et al. The Anticipated
589 Clinical and Economic Impact of 90-90-90 in South Africa. *Annals of internal medicine*.
590 2016;165(5):325–333. Available from: [https://www.ncbi.nlm.nih.gov/pmc/articles/-](https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5012932/pdf/nihms784208.pdf)
591 [PMC5012932/pdf/nihms784208.pdf](https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5012932/pdf/nihms784208.pdf).
- 592 [36] Bershteyn A, Klein DJ, Eckhoff PA. Age-dependent partnering and the HIV transmission
593 chain: a microsimulation analysis. *Journal of the Royal Society, Interface*. 2013
594 nov;10(88):20130613. Available from: [http://rsif.royalsocietypublishing.org/cgi/doi/10.1098/-](http://rsif.royalsocietypublishing.org/cgi/doi/10.1098/rsif.2013.0613)
595 [rsif.2013.0613](http://rsif.royalsocietypublishing.org/cgi/doi/10.1098/rsif.2013.0613).
- 596 [37] McCormick AW, Abuelezam NN, Rhode ER, Hou T, Walensky RP, Pei PP, et al.
597 Development, calibration and performance of an HIV transmission model incorporating natural
598 history and behavioral patterns: application in South Africa. *PloS one*. 2014 may;9(5):e98272.
599 Available from: <http://dx.plos.org/10.1371/journal.pone.0098272>.
- 600 [38] Johnson LF, Kubjane M, Moolla H. MicroCOSM: a model of social and structural drivers of
601 HIV and interventions to reduce HIV incidence in high-risk populations in South Africa. *bioRxiv*.
602 2018;p. 310763.
- 603 [39] Gopalappa C, Farnham PG, Chen YH, Sansom SL. Progression and Transmission of
604 HIV/AIDS (PATH 2.0). *Medical decision making : an international journal of the Society for*
605 *Medical Decision Making*. 2017 feb;37(2):224–233.

- 606 [40] Bakker R, Korenromp E, Meester E, Van Der Ploeg C, Voeten H, Van Vliet C, et al. Stdsim:
607 A microsimulation model for decision support in the control of hiv and other stds. Sexually
608 Transmitted Diseases. 2000;27(10):652.
- 609 [41] Marshall_Labs. Treatment of infectious transmissions through agent-based network.
610 2017;Available from: <https://titan-documentation.readthedocs.io/en/latest/index.html>.
- 611 [42] Bershteyn A, Gerardin J, Bridenbecker D, Lorton CW, Bloedow J, Baker RS, et al.
612 Implementation and applications of EMOD, an individual-based multi-disease modeling
613 platform. Pathogens and disease. 2018;76(5):fty059.
- 614 [43] Penny MA, Galactionova K, Tarantino M, Tanner M, Smith TA. The public health impact of
615 malaria vaccine RTS,S in malaria endemic Africa: Country-specific predictions using 18 month
616 follow-up Phase III data and simulation models. BMC Medicine. 2015 jul;13(1):170.
- 617 [44] Chang ST, Chihota VN, Fielding KL, Grant AD, Houben RM, White RG, et al. Small
618 contribution of gold mines to the ongoing tuberculosis epidemic in South Africa: a modeling-
619 based study. BMC medicine. 2018 apr;16(1):52.
- 620 [45] Fojo AT, Kendall EA, Kasaie P, Shrestha S, Louis TA, Dowdy DW. Mathematical Modeling
621 of Chronic Infectious Diseases: Unpacking the Black Box. In: Open forum infectious
622 diseases. vol. 4. Oxford University Press US; 2017. p. ofx172.
- 623 [46] Gilbert JA, Meyers LA, Galvani AP, Townsend JP. Probabilistic uncertainty analysis of
624 epidemiological modeling to guide public health intervention policy. Epidemics. 2014;6:37–45.
- 625 [47] Caro JJ, Briggs AH, Siebert U, Kuntz KM. Modeling good research practices—overview: a
626 report of the ISPOR-SMDM Modeling Good Research Practices Task Force–1. Medical Decision
627 Making. 2012;32(5):667–677.
- 628 [48] Fehr J, Heiland J, Himpe C, Saak J. Best practices for replicability, reproducibility and
629 reusability of computer-based experiments exemplified by model reduction software. arXiv
630 preprint arXiv:160701191. 2016;.

- 631 [49] Taylor DC, Pawar V, Kruzikas D, Gilmore KE, Pandya A, Iskandar R, et al. Methods of
632 model calibration. *Pharmacoeconomics*. 2010;28(11):995–1000.
- 633 [50] Gerberry DJ. An exact approach to calibrating infectious disease models to surveillance
634 data: The case of HIV and HSV-2. *Mathematical Biosciences & Engineering*. 2018;15(1):153–179.
- 635 [51] Hodges JS. Six (or so) things you can do with a bad model. *Operations Research*.
636 1991;39(3):355–365.
- 637 [52] Kenyon CR, Delva W, Brotman RM. Differential sexual network connectivity offers a
638 parsimonious explanation for population-level variations in the prevalence of bacterial
639 vaginosis: a data-driven, model-supported hypothesis. *BMC women’s health*. 2019;19(1):8.
- 640 [53] Delva W, Leventhal GE, Helleringer S. Connecting the dots: network data and models in
641 HIV epidemiology. *Aids*. 2016;30(13):2009–2020.
- 642 [54] Karnon J, Vanni T. Calibrating models in economic evaluation. *Pharmacoeconomics*.
643 2011;29(1):51–62.
- 644 [55] Kopec JA, Finès P, Manuel DG, Buckeridge DL, Flanagan WM, Oderkirk J, et al. Validation
645 of population-based disease simulation models: a review of concepts and methods. *BMC public*
646 *health*. 2010;10(1):710.
- 647 [56] Talts S, Betancourt M, Simpson D, Vehtari A, Gelman A. Validating Bayesian inference
648 algorithms with simulation-based calibration. *arXiv preprint arXiv:180406788*. 2018;.
- 649 [57] Srikrishnan V, Keller K. How much data are needed to calibrate and test agent-based
650 models? *arXiv preprint arXiv:181108524*. 2018;.
- 651 [58] Zhang H, Vorobeychik Y. Empirically grounded agent-based models of innovation
652 diffusion: a critical review. *Artificial Intelligence Review*. 2019;p. 1–35.
- 653 [59] Murray EJ, Robins JM, Seage III GR, Lodi S, Hyle EP, Reddy KP, et al. Using observational
654 data to calibrate simulation models. *Medical Decision Making*. 2018;38(2):212–224.

655 [60] Lee JS, Filatova T, Ligmann-Zielinska A, Hassani-Mahmooei B, Stonedahl F, Lorscheid I,
656 et al. The complexities of agent-based modeling output analysis. Journal of Artificial Societies
657 and Social Simulation. 2015;18(4):4.

658 [61] Punyacharoensin N, Edmunds WJ, De Angelis D, White RG. Mathematical models for the
659 study of HIV spread and control amongst men who have sex with men. European journal of
660 epidemiology. 2011;26(9):695.

661

662 **Supporting information captions**

663 **S1 Table. Articles included for review**

664 **S2 Table. Description of calibration algorithms**

665 **S1 Text. Obtaining parameter uncertainty using an optimisation algorithm, quoted**
666 **from Sauboin et al.**

667 **S2 Text. Selected quotes of rationales for choosing model calibration method**

668 **S1 Appendix. Parameter search strategies by disease and year of publication**

669 **S2 Appendix. Histograms and plots for counts of targets, calibrated parameters**
670 **and the size of the simulated population**