



Deposited via The University of York.

White Rose Research Online URL for this paper:

<https://eprints.whiterose.ac.uk/id/eprint/104877/>

Version: Accepted Version

Article:

Read, Mark N, Alden, Kieran, Rose, Louis M et al. (2016) Automated multi-objective calibration of biological agent-based simulations. *Interface*. ISSN: 1742-5662

<https://doi.org/10.1098/rsif.2016.0543>

Reuse

Items deposited in White Rose Research Online are protected by copyright, with all rights reserved unless indicated otherwise. They may be downloaded and/or printed for private study, or other acts as permitted by national copyright laws. The publisher or other rights holders may allow further reproduction and re-use of the full text version. This is indicated by the licence information on the White Rose Research Online record for the item.

Takedown

If you consider content in White Rose Research Online to be in breach of UK law, please notify us by emailing eprints@whiterose.ac.uk including the URL of the record and the reason for the withdrawal request.

Automated Multi-Objective Calibration of Biological Agent-Based Simulations

Mark N. Read^{1,2,*}, Kieran Alden³, Louis M. Rose⁴, Jon Timmis³

¹School of Life and Environmental Sciences, The University of Sydney, Australia.

²Charles Perkins Centre, The University of Sydney, Australia.

³Department of Electronics, The University of York, UK.

⁴Department of Computer Science, The University of York, UK.

*Corresponding author: mark.read@sydney.edu.au

Keywords: Computational Biology — Agent-Based Simulation — Calibration — Multi-Objective Optimization — Experimental Autoimmune Encephalomyelitis

Abstract

Computational agent-based simulation is increasingly used to complement laboratory techniques in advancing our understanding of biological systems. Calibration, the identification of parameter values that align simulation with biological behaviours, becomes challenging as increasingly complex biological domains are simulated. Complex domains cannot be characterised by single metrics alone, rendering simulation calibration a fundamentally multi-metric optimisation problem that typical calibration techniques cannot handle. Yet calibration is an essential activity in simulation-based science; the baseline calibration forms a control for subsequent experimentation, and hence is fundamental in the interpretation of results. Here we develop and showcase a method, built around multi-objective optimisation, for calibrating agent-based simulations against complex target behaviours requiring several metrics (termed *objectives*) to characterise. Multi-objective calibration delivers those sets of parameter values representing optimal tradeoffs in

simulation performance against each metric, in the form of a Pareto front. We use MOC to calibrate a well-understood immunological simulation against both established *a priori* and previously unestablished target behaviours. Further, we show that simulation-borne conclusions are broadly, but not entirely, robust to adopting baseline parameter values from different extremes of the Pareto front, highlighting the importance of MOC’s identification of numerous calibration solutions. We devise a method for detecting overfitting in a multi-objective context, not previously possible, used to save computational effort by terminating MOC when no improved solutions will be found. MOC can significantly impact biological simulation, adding rigour to and speeding up an otherwise time-consuming calibration process, and highlighting inappropriate biological capture by simulations that cannot be well calibrated. As such, it produces more accurate simulations that generate more informative biological predictions.

1 Introduction

Computational modelling and simulation has emerged as a tool for investigating a wide range of biological systems, spanning immunology [1][2], drug and intervention design [3][4], developmental biology [5], and ecology [6]. Biological simulation is particularly insightful when used in complement with traditional methods, such as wet-lab *in vivo* and *in vitro* work; laboratory work generates experimental data and suggests hypotheses that can be evaluated by way of their integration with simulation, which in turn can suggest further experiments or highlight areas of lacking knowledge [7][8]. Well designed, biologically-accurate simulations provide detailed spatio-temporal insight, facilitating observations and assays not possible in the real system; simulation experiments are unhampered by the ethical, practical and financial considerations inherent in biological experimentation. Research programs integrating wet-lab and simulation methods can offer a greater return on animal experimentation by generating additional insight, and hence easing the burden on experimental animals, in line with the ‘3Rs’ principles (Replacement, Reduction and Refinement).

The agent-based simulation (ABS) paradigm permits detailed and nuanced simulation of biological systems [9][3]. Simulation components are represented as explicit individual entities, *agents*, with unique states that exist within a spatial environment. Rules

59 specifying agent dynamics and the consequences of interaction are provided, and simu-
60 lation execution allows the system-level consequences of agent-level manipulations to be
61 observed. ABS incorporates stochastic events, and therein reflects the heterogeneity of
62 real world natural systems. There is scope for specifying very detailed interactions using
63 ABS, at the expense of generating large numbers of parameters: 50+ is not uncommon.

64 Drawing biologically meaningful conclusions from simulation requires that the map-
65 ping of the simulation to the biology is known. This can prove problematic for two
66 reasons. First, simulations are abstract representations of their corresponding real world
67 systems. For example, there exist at least 19 varieties of T cell, a vital component of the
68 immune system [10]. However, rather than fully capture all their nuanced differences, a
69 simulation is more likely to represent an abstracted subset thereof. As such, experimen-
70 tal measurements on a real world T cell cannot be assumed to translate directly to its
71 simulation counterpart. Second, complex biological systems are the subject of simulation
72 precisely because they are incompletely understood, meaning that the real world data sup-
73 porting simulation design decisions and corresponding parameter values may not exist.
74 Calibration is a critical activity in establishing the link between simulation and biology;
75 parameter values that align simulation and real-world dynamics are identified. Further-
76 more, an inability to provide a good alignment points to simulation design that does not
77 appropriately capture the biology. Calibration is used to establish a baseline simulation
78 dynamic used as a control in subsequent experimentation, and finding appropriate values
79 is important. Different parameter values will yield different simulation dynamics, and as
80 such influence the conclusions drawn from experiments.

81 A number of approaches to calibration exist, including manual calibration [11], evo-
82 lutionary algorithms [12][13], maximum likelihood estimation and various forms of re-
83 gression [14]. These techniques identify parameter values by employing a single metric
84 to align simulation dynamics with those of the real world system. However, complex
85 biological system dynamics are not well characterised by single metrics alone. They con-
86 stitute many different types of interacting component, and encompass both positive and
87 negative feedbacks. They are highly redundant: a single component can perform many
88 functions and any one function can be performed by several components [15][16]. As such,
89 calibration of a complex system simulation is fundamentally a multi-metric optimisation
90 problem; several metrics of a simulation's alignment with the biology must be simulta-

91 neously considered when evaluating putative parameter values. Consider, for example,
92 cellular motility, which underlies many biological processes arising from cellular interac-
93 tion. Which targets a given cell interacts with depends on both its speed and directional
94 persistence; accurately modelling this process requires that metrics of both be considered.

95 In this paper, we position multi-objective optimisation-based calibration (MOC: multi-
96 objective calibration) as an important enabling technology for simulation-based biological
97 investigation. Given its abstractive nature, a simulation undergoing calibration will not
98 perfectly replicate all aspects of the biology. As such, putative simulation parameter value
99 sets will exhibit tradeoffs in their reproduction of aspects of the biology, excelling in some
100 at the expense of others. In this context, a metric quantifying a simulation’s capture
101 of a specific aspect of the biology is termed an *objective*. Through the use of Pareto
102 fronts (defined in Section 3), MOC explicitly tracks the collection of simulation parameter
103 sets exhibiting optimal tradeoffs between objectives. It is unknown if adopting baseline
104 parameter values from different regions of the Pareto front will deliver fundamentally
105 different conclusions from simulation-based experiments. The answer to this question is
106 likely problem-specific, and the use of MOC allows this issue to be addressed by exposing
107 a full range of Pareto-equivalent solutions.

108 Here we investigate multi-objective optimisation, specifically the NSGA-II algorithm [17],
109 in calibrating an established immunological simulation: ARTIMMUS [18]. ARTIMMUS
110 simulates Experimental Autoimmune Encephalomyelitis (EAE), a mouse model of multi-
111 ple sclerosis [19][20]. It is a complex simulation, encompassing seven distinct cell popula-
112 tions that interact across five organs, and constituting 72 parameters. Its successful prior
113 manual calibration renders it an effective test case for evaluating MOC’s applicability to
114 simulation calibration. We demonstrate the successful calibration of ARTIMMUS using
115 five objectives (Section 4): a range of solutions to the calibration problem, offering optimal
116 tradeoffs against calibration objectives, are generated. Furthermore, we demonstrate that
117 conclusions drawn from a simulation-based experiment can vary depending on exactly
118 which calibration solution is adopted (Section 5). Hence, different calibration solution
119 parameter values can vary downstream conclusions, highlighting MOC’s value in making
120 these multiple solutions explicit. We show that MOC is equally applicable in generating
121 simulation initial condition values: cellular population sizes as simulation launch. We
122 proceed to demonstrate that MOC can identify parameter and initial condition values

123 that deliver previously unknown simulation dynamics, highlighting its potential beyond
124 this well understood test case (Section 6). Lastly, we consider strategies for formulat-
125 ing stopping criteria for MOC, thereby preventing over-fitting and wasted computational
126 expense when apparent improvements in simulation calibration are likely due to stochas-
127 tic sampling rather than genuinely superior parameter values (Section 7). We begin by
128 introducing ARTIMMUS (Section 2), and the MOC methodology (Section 3).

129 **2 A testbed for calibrating biological simulations**

130 ARTIMMUS is an agent-based simulation of an EAE protocol wherein mice induced into
131 autoimmunity undergo a natural recovery from disease, and are thereafter resistant to
132 disease re-induction [18][21][22]. ARTIMMUS was created, in part, to further probe the
133 cellular interactions mediating this recovery [23][24]. It has been used to explore the
134 mechanisms through which splenectomy, the removal of the spleen, a primary immune
135 organ, exacerbates disease severity, and predict the outcome of T cell interaction-blocking
136 drugs [18]. It was conceived through a collaboration of immunologists and computer
137 scientists, and developed through a principled approach focusing on documenting how
138 biological concepts are translated into computer code: the CoSMoS process [32]. It is
139 written in the Java programming language.

140 ARTIMMUS has previously undergone a by-hand, manual calibration [11], and was
141 shown to reflect the dynamics of the real world disease [18]. The process demanded close
142 collaboration between the simulation developer and an immunologist who informed the
143 work, helping bridge biological data and concepts to simulation constructs and output.
144 This manual calibration took two weeks, and entailed an iterative process through which
145 simulation code and parameter value changes that might explain perceived discrepancies
146 between simulation and biological system dynamics were identified and explored in turn.
147 Those best aligning simulation with biological dynamics were adopted before repeating
148 the process. This calibration approach is akin to a non-population, manual, greedy local
149 search wherein the best immediate improvement is always adopted.

150 Despite delivering a well-calibrated result for ARTIMMUS, this calibration search
151 strategy presents several potential pitfalls. It is entirely plausible that the manual search
152 does not find the global optimum parameter set that best aligns simulation dynamics with

153 those of the biological system. As a greedy search strategy, its result is highly dependent
154 on the search's starting position, and complex landscapes where one parameter's influence
155 on simulation dynamics critically depends on the values held by others are particularly
156 challenging. The existence of multiple solutions to the calibration problem can go entirely
157 undetected. Lastly, manual calibration is time consuming, and agent-based simulation's
158 stochastic nature furthers compound these challenges. It is these issues that collectively
159 motivated the present automated MOC approach.

160 Here we provide a brief summary of EAE and ARTIMMUS to aid understanding of
161 the sections that follow; a comprehensive description may be found in the supplementary
162 materials of [18]. Figure 1A provides an abstract overview of the major cell types in-
163 volved in EAE, and their relationships to one another. EAE is induced through injection
164 of neuronal fragments which are internalized by dendritic cells (DCs) which then direct
165 the growth of a T cell population (CD4Th1, abbreviated to Th1) targeting these frag-
166 ments. These Th1 cells enter the central nervous system (CNS), where they stimulate
167 CNS-resident macrophages into secreting $\text{TNF-}\alpha$, which in turn damages neurons. The
168 resultant neuronal fragments are internalized by further populations of DCs, which direct
169 further Th1 activities, perpetuating the autoimmune cycle. Recovery from autoimmunity
170 is through the actions of two populations of regulatory T cell, CD4Treg and CD8Treg
171 cells, so named as they regulate the activities of other T cells. The natural life-cycle of a
172 Th1 cell results in its eventual death and internalization by DCs, which derive fragments
173 therefrom and direct the growth of CD4Treg and CD8Treg cells targeting the Th1 cell
174 population. CD4Tregs play an essential role in facilitating the development of CD8Treg
175 cells. CD8Treg cells can directly kill Th1 cells, interrupting their natural life-cycle and
176 preventing the perpetuation of autoimmunity. Th2 cells directly compete with Th1 cells,
177 as both arise from a common progenitor and they each perform downstream activities
178 that promote their own development. The reduced severity of the autoimmune environ-
179 ment arising from the action of CD8Treg cells favours the growth of Th2 cells over Th1
180 cells, which do not directly harm neurons and hence do not contribute to this autoimmune
181 process. Figure 1B shows a time-series graph of T cell population sizes in ARTIMMUS.

182 **3 Multi-objective calibration (MOC) methodology**

183 We present here an overview of the multi-objective calibration (MOC) concept, detailing
184 how we employ multi-objective optimisation technology to calibrate simulation parameters
185 and initial conditions. A graphical overview is supplied in Figure 2.

186 Firstly, we define the desired (*target*) ARTIMMUS dynamics, Figure 2A. In this
187 manuscript targets are expressed as peak cell population sizes, the times at which those
188 peaks occur, or the cell population sizes at a given time. Target dynamics might repre-
189 sent known biological results to be reproduced, or hypothetical outcomes of interest. In
190 this study we adopt the dynamics of a previous manual calibration of ARTIMMUS, so
191 as to evaluate MOC on a well-understood problem; thereafter we employ MOC to obtain
192 hypothetical dynamics not known possible *a priori*. We note that many other aspects
193 of simulation performance can constitute target dynamics, depending on the context and
194 simulation being calibrated. The expression of targets as distributions reflects the stochas-
195 tic nature of biological systems and agent-based simulations, wherein repeat experiments
196 can yield slightly different results.

197 MOC seeks to identify parameter values that best align simulation with target dynam-
198 ics. As such, we define metrics, termed *objectives*, that quantify the alignment between the
199 two. As illustrated in Figure 2B (left), we employ the the non-parametric Kolmogorov-
200 Smirnov statistic in our objectives, which quantifies the difference in target and simulation
201 dynamics for a given set of simulation parameter values. Rather than contrasting the me-
202 dians of two distributions, as many statistics do, the KS statistic quantifies the biggest
203 distance between two distributions' cumulative distribution functions. As such, its use
204 here facilitates the calibration of a distribution's shape, not simply its median or mean.
205 We consider this a strength of our approach; as may be seen in the sections that follow,
206 MOC is capable of reproducing distributions of behaviour, not simply averages. Each
207 set of simulation parameter values is termed a 'candidate solution', and its corresponding
208 simulation performance is evaluated against each objective individually. By evaluating
209 many candidate solutions we identify regions of parameter space providing close align-
210 ment with target dynamics (Figure 2B, right). Importantly, the regions that satisfy each
211 objective differ. In practice, it is computationally intractable to fully explore parameter
212 space as suggested by the heatmaps in this Figure, particularly when many parameters
213 are investigated. Instead, a heuristic (guided) search strategy is employed that samples

214 parameter space, evaluates performance, and decides from where to extract the next can-
215 didate solutions based on the results. In this study we employ NSGA-II as our guided
216 search engine [17], but we believe other multi-objective optimisation technologies could
217 be successfully substituted. NSGA-II maintains a population of candidate solutions, and
218 employs (heavily abstracted) principles of genetic recombination, mutation and natural
219 selection to generate and evaluate successive generations of superior candidate solutions.
220 Hence, NSGA-II is an iterative algorithm. We refer the readers to [17] for more detail on
221 NSGA-II. Here we have employed the ‘inspyred’ python module NSGA-II implementation.

222 We identify those candidate solutions that constitute optimal tradeoffs in performance
223 against each objective, referred to simply as *solutions*, Figure 2C. The set of solutions
224 is termed the *Pareto front*. These solutions are *Pareto-equivalent*: no solution has been
225 found that offers an improvement in one objective without a worsening in another. Pareto-
226 equivalent solutions may reside in disparate regions of parameter space, and the ability to
227 recognize this is a key strength of MOC. Though these regions of parameter space may be
228 Pareto-equivalent for the given target simulation behaviour, they could yield very different
229 behaviours when subjected to further downstream experimentation, and as such lead to
230 different simulation-borne conclusions. In this study we investigate this phenomenon for
231 a given experiment in ARTIMMUS.

232 We note that it is possible to derive a great many targets and objectives for complex
233 system simulations. Increasing the number of objectives increases the difficulty of the
234 calibration problem, and the computational resource required to address it; in the field
235 of optimisation this is known as the ‘curse of dimensionality’. Hence, employing fewer,
236 uncorrelated objectives is considered good practice: it encourages the identification of
237 good quality solutions whilst minimising the resources required to do so.

238 **3.1 Selecting candidates from the Pareto front**

239 Upon completion MOC delivers a Pareto front of Pareto-equivalent solutions, representing
240 optimal tradeoffs between the calibration objectives. Deciding which solution adopt as the
241 baseline simulation parameter values is an application-specific problem. For the present
242 study we have developed a function, $\Lambda(c)$, which assesses candidate solution c against
243 the criteria below. We select the candidate with the lowest Λ value when presenting the
244 results of calibration below. Λ is calculated as follows.

245 Let Ω represent the set of calibration objectives, and $KS_o(c)$, $o \in \Omega$ as the correspond-
 246 ing Kolmogorov-Smirnov score for candidate c on objective o . $\overline{KS}(c)$ represents the mean
 247 objective score for candidate c . The Λ score is calculated as:

$$\Lambda(c) = \alpha \cdot \overline{KS}(c)^2 + \sum_{o \in \Omega} (KS_o(c) - \overline{KS}(c))^2 \quad (1)$$

248 Low Λ scores are achieved through low mean objective KS scores, and balanced KS
 249 scores across all objectives. α specifies the relative importance of these two components.
 250 When $\alpha = 1$, both measures contribute equally to Λ . Lower mean KS scores are prioritised
 251 with $\alpha > 1$, and *vice versa*. We employ $\alpha = 1$ throughout. We note that Λ is unit-less, and
 252 as such is not explicitly reported here; it is used only to extract one candidate solution
 253 from a Pareto front, presented as the chief result of calibration in the results that follow.

254 4 Successful re-calibration of ARTIMMUS

255 We demonstrate MOC by re-calibrating ARTIMMUS, taking as target dynamics those of
 256 the previous manually-calibrated simulation dynamics [18]. As these dynamics are known
 257 to be obtainable, and at least one set of parameter values that produce them are known,
 258 we are able to evaluate MOC's performance.

259 With 5 objectives MOC successfully reproduced the manually-calibrated ARTIMMUS
 260 dynamics, as demonstrated in Figure 3. The objectives used were:

- 261 • the peak Th1 cell population size (Figure 3B)
- 262 • the time at which the peak occurred (Figure 3C)
- 263 • the Th2 population size at 30 days (Figure 3D)
- 264 • the peak population sizes of both CD4Treg and CD8Treg cells (Figures 3E and F).

265 The corresponding target distributions of values are also shown in Figure 3.

266 Each candidate solution generated by NSGA-II was assessed through 200 replicate
 267 simulation executions. The target distributions against which candidates are contrasted
 268 are derived from 500 replicates generated with the previous manual-calibration parameter
 269 values. The manual-calibration's replicates need be executed once only and stored, they
 270 do not change. In contrast, assessment of candidates is computationally costly because

271 so many are generated; a figure of 200 replicates per candidate was selected to strike a
272 balance between experimental sensitivity and computational cost. A previous analysis of
273 parametric perturbation in ARTIMMUS established that contrasting distributions com-
274 prising 200 replicate executions was sufficient to detect ‘small’ changes in $\frac{2}{3}$ of simulation
275 behaviour metrics, and ‘medium’ in the remainder [11]. Hence, we consider 200 replicates
276 to offer sufficient sensitivity in differentiating candidate performances. These effect size
277 categories arise from the analysis’s use of the Vargha-Delaney *A test* [25], which provides
278 interpretation guidelines. For reference, the *A test* is a non-parametric effect magnitude
279 test representing the probability that a randomly selected member of one distribution is
280 larger than a randomly selected member of the other. An *A test* score of 0.5 indicates
281 the two distributions are indistinguishable (using this test). Values of 1 and 0 indicate
282 no overlap in the two distributions. A single calibration exercise required around 5 days
283 on a dedicated computational cluster able to execute 120 simulations simultaneously;
284 each single simulation replicate takes around 2-10 minutes to execute, depending on the
285 parameter values used.

286 We have successfully applied MOC to both ARTIMMUS parameter values and ini-
287 tial conditions, but focus here on the former. Initial condition calibration results are
288 reported in the supplementary materials. Calibration was performed over 8 ARTIMMUS
289 parameters which all pertain to presentation of substances to T cells, particularly Th1
290 and Th2 cells, and their resultant development. The biology captured in these param-
291 eters is outlined in supplementary Figure S1, and we note that a through understanding
292 of this biology is not required to appreciate our results. These parameters were selected
293 for the reasons that ascertaining their values experimentally would be challenging and
294 they all relate to a critical aspect of the biology: the perpetuation of autoimmunity, and
295 (for some) it’s amelioration (as Treg cell development is also directed by DCs). Hence,
296 by successfully calibrating parameter values that are highly influential on simulation dy-
297 namics we demonstrate MOC’s potential. Parameters were given a constrained range of
298 values that the MOC process could assign, being zero to twice their manually-calibrated
299 range, as shown in Table 1. In exploring the space of putative parameter values, NSGA-II
300 maintained a population of 64 candidate solutions which were subject to genetic recombi-
301 nation and mutation (see [17]) over 32 generations of natural selection, wherein only the
302 best 64 solutions (i.e. those on or near the Pareto front) were retained in the successive

303 generation.

304 This calibration exercise was repeated three times for both parameters and initial
305 conditions. Figure 3 shows the solution with the lowest Λ score from one such parameter
306 calibration. The remaining two are shown in supplementary Figures S2 and S3. The
307 calibrated simulation dynamics closely resemble the target distributions in all cases. The
308 three parameter calibration exercises generated, respectively, Pareto fronts constituting
309 82, 87 and 112 Pareto-equivalent solutions. The ranges of parameter values represented
310 across the Pareto fronts' solutions in each independent calibration exercise are shown in
311 Figure 4, as are the baseline manually-calibrated values. In all but one case the baseline
312 parameter value sat within the range of non-outlier MOC-derived values, the exception
313 being *Th1_diff80* in exercise 3. Hence, we conclude that MOC is an effective means
314 of calibration: it has repeatedly reproduced ARTIMMUS dynamics that were known
315 possible, and has identified similar solutions, in the form of parameter values, that do so.

316 Next we investigated how the space of ARTIMMUS parameter values relates to the
317 space of successful target dynamic reproductions, i.e., tradeoffs in objective values. We
318 find statistically significant ($p < 0.01$) differences between calibration exercises' distribu-
319 tions of calibrated parameter values for 7 of 8 parameters, Figure 4. This corresponds
320 to 19 of 24 (79%) of pairwise comparisons. Further, 75% (18/24) pairwise comparisons
321 register a KS value ≥ 0.3 . For context, a KS value of 1.0 indicates no overlap between 2
322 distributions. In contrast, this degree of variation is not observed in Pareto fronts' objec-
323 tive values, depicted in Figure 5. Here we instead find statistically significant differences
324 in only 27% (4/15) pair-wise calibration comparisons, and only 27% (4/15) of compar-
325 isons register $KS \geq 0.3$. We find no evidence of objectives that are harder to calibrate
326 than others; the smallest objective values are < 0.05 in all cases, and the median objective
327 values all lie under 0.17.

328 Together, these data suggest a redundancy in the ability for parameter values to
329 deliver particular objective scores. This corresponds to a landscape wherein parameter
330 values mapped to objective values is relatively flat, as a wide range of ARTIMMUS
331 parameter values deliver relatively similar objective scores. The results of using MOC
332 to calibrate ARTIMMUS initial conditions are reported in supplementary Section S1, and
333 supplementary Figures S4, S5 and S6. They are qualitatively identical to our findings in
334 calibrating parameters, and support the conclusions drawn here.

335 An obvious question is, why does MOC not deliver any perfectly calibrated solutions,
336 wherein all objective scores are 0.0? The best solutions, determined by their minimal Λ
337 values, in each calibration exercise are shown in Table 2. Objective KS values ranged
338 from 0.05 to 0.14 (and 0.03 to 0.12 for initial conditions). We attribute the inability to
339 deliver a perfect calibration to the stochastic nature of ARTIMMUS, wherein 200 replicate
340 executions for a given candidate yields sufficient variation so as to deliver objective KS
341 scores of ≥ 0.05 . There is a risk that improvements in objective KS values that are already
342 so small cannot be confidently attributed to an actual improved simulation calibration, as
343 opposed to stochastic variation between simulation replicates. Section 7, below, explores
344 a method for terminating the MOC process on the premise that further effort will not
345 deliver better quality solutions.

346 These data collectively highlight the challenges in exactly calibrating (i.e. $KS=0.0$)
347 simulations to several objectives simultaneously. As such, we consider in the next Sec-
348 tion the implications on experimental results of adopting baseline simulation values from
349 different extremes of the Pareto front.

350 5 Scientific significance of imperfect calibration

351 As demonstrated above, MOC delivers a host of solutions to a given calibration problem,
352 each representing an optimal tradeoff in calibration criteria (see Figure 2). It falls on the
353 simulation developer to decide which to adopt baseline parameter values in subsequent
354 experimentation. There is a risk that whilst calibration solutions lying in different regions
355 of parameter space give rise to Pareto-equivalent solutions, they do not behave in a
356 consistent manner when further experiments are performed. In such a case, a simulation-
357 based experiment would lead to different conclusions depending on which calibration
358 result was adopted as the baseline. In this section we investigate the extent to which this
359 phenomenon holds.

360 The manually-calibrated ARTIMMUS simulation was previously used to elucidate the
361 effect of removing a central immune organ, the spleen (a *splenectomy*), in EAE-induced
362 animals [18]. Previous experiments had demonstrated that splenectomy in rats prior to
363 the induction of EAE increased the mortality rate and hampered recovery [26]. Simulating
364 splenectomy in ARTIMMUS revealed the spleen as a primary site for the generation of

365 autoimmunity-combating CD4Treg and CD8Treg cells. The reduced Treg populations
366 resulting from the spleen’s removal prior to EAE-induction were unable to completely
367 abrogate the autoimmunity-inducing Th1 populations, allowing for their re-expansion,
368 and thus facilitating increased disease severity and relapses.

369 Here we explore whether the results of splenectomy in ARTIMMUS differ when base-
370 line parameter values are adopted from disparate extremes of the Pareto front. The
371 experimental procedure is highlighted in Figure 6. First, Pareto front solutions represent-
372 ing the extreme values, both low and high, of objective KS measures are identified. These
373 solutions represent extremes in the range of simulation dynamics encapsulated within the
374 Pareto front. For each solution 200 simulation replicates are performed for both control
375 and splenectomy groups. Key performance indicators (KPI) are extracted from the resul-
376 tant distributions of 200 simulation executions in each group. The performance indicators
377 used are identical to those of the original ARTIMMUS splenectomy experiment [18]: the
378 peak population sizes for each T cell population in the simulation, the times at which
379 these peaks are reached, and the number of Th1 cells remaining at day 40 (giving a total
380 of 9). For each KPI, the distributions of values obtained for control and splenectomy
381 groups are contrasted using the Vargha-Delaney *A test* [25], as per the original experi-
382 ment [18]. This procedure is repeated for each of the three calibration exercises reported
383 in Section 4. The resultant *A test* scores are shown in Figure 6’s tables. Also shown, for
384 context, are the *A test* scores of the original ARTIMMUS experiment [18].

385 Broadly speaking, the splenectomy results generated by Pareto-equivalent solutions are
386 consistent with one another, and with the original experiment. There exceptions, however,
387 wherein differences in *A test* scores reported for solution and the original experiment
388 differed substantially: g23c60 in exercise 1, and g6c35 and g30c58 in exercise 2. These
389 differences occurred for ‘Th1 at 40d’, ‘Th2 peak’ and ‘Th2 Time’ KPIs. Of interest, three
390 of these solutions were obtained from the region of the Pareto front where alignment with
391 target Th2 peak population size was poorest. In the case of g23c60 and g6c35, exercises 1
392 and 2 respectively, the parameter values were sufficient to return Th1 population size at
393 40 days to control group levels, despite the splenectomy ($A=0.58$ and 0.56 ; 0.5 indicates
394 no difference). This is significant, as the principle conclusion of the original experiment
395 was that splenectomy reduces Treg population sizes to levels unable to suppress Th1 cell
396 populations and abrogate autoimmunity. The time series T cell population dynamics of

397 both these solutions under control and splenectomy are shown in supplementary Figure S9.
398 In both cases the peak Th1 population sizes are smaller than in the original experiment
399 (see Figure 6)), and the Th2 population sizes are substantially larger. Based on this we
400 hypothesize that despite reduced Treg population sizes resulting from splenectomy, the
401 altered balance between Th1 and Th2 populations which compete with one another is
402 sufficient to abrogate the Th1 population at day 30 in these solutions.

403 Supporting the notion that solutions' results are relatively consistent, the direction of
404 change in solutions' KPIs resulting from splenectomy differs from the original experiment
405 in only a minority of cases. Further, this occurs only in KPIs for which the original
406 experiment reports a comparatively small change between splenectomy and control, the
407 largest being in exercise 2 when the original experiment reports a change of $A=0.66$, which
408 was not interpreted as significant.

409 We have conducted the same investigation on Pareto-equivalent solutions generated
410 under the three independent initial condition calibration exercises (supplementary Section
411 S1). Detailed analysis is reported in supplementary Section S2 and Figure S10; briefly,
412 divergences between initial condition solution and original experiments were smaller than
413 reported here for parameters. We take this to indicate that the initial parameters in-
414 vestigated were less influential on simulation behaviour than the parameters investigated
415 here.

416 In summary, the conclusions that would be drawn from adopting baseline parameters
417 values from disparate Pareto-equivalent solutions are mostly, but not completely, consis-
418 tent with one another and with the original splenectomy experiment. There were two
419 notable exceptions, and they underscore the importance of considering the range of sim-
420 ulation performances that satisfy a calibration exercise. Making these explicit through
421 Pareto fronts is a strength of the MOC approach. It remains important to, where possi-
422 ble, further evaluate Pareto-equivalent solutions in the context of domain knowledge and
423 expertise, which might have ruled out the two exceptions noted above, as the Th2 popula-
424 tion size is abnormally large compared to the Th1 population. Where this is not possible,
425 where no grounds to discard some Pareto-equivalent solutions exist, we advise that ex-
426 periments are performed in replicate adopting a wide range of calibration solutions and
427 that conclusions are drawn after taking stock of the full range of results generated. This
428 is particularly important if quantitative, rather than qualitative, results are sought; our

429 present data show more divergence between calibration solutions and original experiment
430 in the quantitative case.

431 **6 Multi-objective calibration delivers previously un-** 432 **seen disease phenotypes**

433 In Section 4, above, MOC successfully reproduced simulation dynamics known to exist
434 by virtue of a prior manual calibration. To further demonstrate MOC's generality and
435 utility, we now derive simulation dynamics not known to exist *a priori*.

436 ARTIMMUS's baseline behaviour constitutes a period of autoimmunity followed by
437 recovery, reflecting typical biological disease [21][22]. However, disease susceptibility
438 and severity vary considerably between mouse strains and between mice within a given
439 strain [27][28]. Furthermore, depletion or incapacitation of CD4Treg and CD8Treg cells
440 leads to exacerbated disease symptoms [29][30]. Here we investigate the capacity for
441 ARTIMMUS to reproduce persisting disease symptoms of varying severity. To reflect
442 potential genetic differences between mouse strains, we calibrate over initial conditions
443 specifying cell population sizes, and a parameter controlling the efficiency of Th1 killing
444 by CD8Treg cells; together comprising 9 variables. In this experiment we are implicitly
445 investigating whether variation in these basal population sizes and the efficiency of the
446 CD8Treg-Th1 killing pathways could explain the differences in autoimmune phenotypes
447 observed between mouse strains and individuals therein.

448 Three persisting disease severities are investigated, ranging from mild to severe. These
449 are captured by defining the distribution of Th1 cells remaining at 60 days as a target for
450 calibration, captured as a Gaussian distribution. Mild, moderate and severe disease are
451 represented with mean (μ) and standard deviation (σ) values of $\mu=50$ & $\sigma=10$, $\mu=200$ &
452 $\sigma=100$, and $\mu=500$ & $\sigma=200$ respectively. To ensure an aggressive onset of autoimmunity,
453 consistent with animal models, a second calibration target distribution of $\mu=1000$ &
454 $\sigma=200$ Th1 cells at 15 days is employed.

455 Each persisting autoimmunity severity is independently calibrated three times, rep-
456 resentatives of which are shown in Figure 7 (the remainder are shown in supplementary
457 Figures S11, S12 and S13). Automated calibration successfully delivers the required me-
458 dian number of cells in most cases, with $KS \leq 0.2$ in 6 of the 9 calibrations. However, the

459 spread of the ‘Th1 cells at 60 days’ distribution for mild persisting disease is notably less
460 well calibrated, with all three calibrations delivering $KS > 0.3$.

461 Together, these data support the general applicability of MOC to problems where a
462 simulation’s ability to deliver a desired dynamic is not known *a priori*. These data also
463 suggest that the heterogeneity in disease severities observed in experimental animals could
464 be attributed to differences in basal population sizes and regulatory pathway efficiency.

465 7 When to stop MOC

466 A key consideration in any optimisation task is the stopping criteria. For MOC, under-
467 pinned by the NSGA-II optimisation algorithm, this equates to determining when to stop
468 calibration.

469 *Overfitting* describes the case where the simulation being calibrated starts to capture
470 the noise in the target distributions, rather than the trends those distributions represent.
471 This is a particular issue when target distributions do not contain many samples, as
472 might be the case if they represent biological experiments (Figure 8A). For example,
473 studies involving experimental animals can require their sacrifice to collect data. As such,
474 it is considered unethical (and is practically cumbersome) to collect hundreds of samples,
475 and 5 to 10 are more typical. These smaller sample sizes are unlikely to perfectly capture
476 the underlying distribution that would emerge if thousands of samples were available.
477 Overfitting is said to have occurred when the calibrated simulation better reflects these
478 5-10 samples than their underlying distribution, as illustrated in Figure 8B.

479 A common strategy in single-objective (not MOC, which is multi-objective) problems
480 for determining when to terminate an optimisation process is to segregate the available
481 data into two parts, termed ‘training’ and ‘validation’ datasets. The training dataset
482 is used as normal to search for improved solutions, akin to MOC’s target data. The
483 validation dataset is used as an independent check for overfitting of solutions to the
484 training data set. Such a case of overfitting is depicted in Figure 8C. Both the training
485 and validation data roughly reflect the underlying distribution, from which they were
486 sampled. The candidate solution more closely resembles the training dataset than either
487 the underlying distribution or the validation dataset, hence, it is overfitted. As illustrated
488 in Figure 8D, in the earlier stages of optimisation successive candidate solutions that

489 better capture the training dataset will also better capture the validation data. It is only
 490 when overfitting starts to occur that performance against the validation data worsens
 491 whilst performance against training data continues to improve. It is at this point that
 492 the optimisation process is best terminated.

493 MOC is, however, a multi-objective optimisation problem, and it is unclear in the lit-
 494 erature how this overfitting detection strategy ought be applied. We propose here a novel
 495 strategy for detecting overfitting in mutli-objective problems based on co-membership of
 496 solutions to both training and validation dataset Pareto fronts (P_t and P_v), maintained
 497 throughout the calibration process (Figure 8E). The overfittedness at a given point in the
 498 optimisation process is reflected in the proportion of P_t members that are not members
 499 of P_v . The following algorithm performs the calculation:

```

500    $m \leftarrow 0$ 
501   for all  $i \in P_t$  do
502     if  $i \in P_v$  then
503        $m \leftarrow m + 1$ 
504     end if
505   end for
506   return  $1 - (m/\text{size}(P_t))$ 

```

507 A proportion of 0 indicates that all training dataset Pareto solutions are also members
 508 of the validation Pareto front. At the other extreme, a value of 1 indicates that the training
 509 dataset Pareto front has been completely over-fitted, as none of its members are Pareto
 510 optimal with respect to the validation dataset. A threshold level of over-fitting at which
 511 the optimisation process (i.e., MOC) is to be terminated can be selected by the simulation
 512 experimenter.

513 We investigated different overfitting thresholds for MOC termination in the three
 514 ARTIMMUS parameter recalibration exercises reported in Section 4 above. An additional
 515 214 simulation replicates using manually-calibrated parameter values were acquired to use
 516 as a validation dataset, constituting a 70-30 (500-215) training-validation data split. The
 517 validation dataset Pareto front for each iteration of the MOC algorithm (generation) was
 518 determined, and the overfittedness calculated. Figure 9A shows how, as MOC progresses,
 519 the proportion of overfitted candidate solutions on the training (target) dataset increases
 520 for each of the three calibration exercises. Figure 9B shows the point at which MOC

521 calibration would have been terminated should a given overfittedness threshold have been
522 selected. Had we employed a overfittedness termination threshold of 0.5, wherein half
523 of the training dataset Pareto front is overfitted, calibration would have terminated at
524 generation 14, 15 or 23 (for exercises 1, 2 and 3 respectively) instead of 32. Given that
525 each of these calibration exercises required around 7 days to complete on a dedicated
526 computing cluster, this speed-up is substantial. We note that these combined training
527 and validation datasets constitute 714 data points, considerably exceeding what might
528 be obtained from real biological experiments. We anticipate that with fewer data points
529 overfitting will occur sooner in the MOC process.

530 8 Discussion

531 Simulation represents a powerful tool to advance the investigation of biological systems,
532 particularly when used in tandem with traditional approaches. As more complex biolog-
533 ical systems become the subject of simulation a challenge in their calibration emerges:
534 complex biological systems cannot be characterised by single metrics alone. There exist
535 technologies capable of identifying parameter values that align simulation dynamics with
536 some desired target, but these operate on single metrics. Even in cases where param-
537 eter values can be ascertained experimentally, seemingly avoiding the need for calibra-
538 tion, the abstract nature of simulation can complicate their direct adoption. Here we
539 have demonstrated how biological agent-based simulation parameter values can be de-
540 rived using multi-objective optimisation, an approach we have termed Multi-Objective
541 Calibration (MOC). Multi-objective optimisation algorithms find solutions to problems
542 simultaneously described by more than one metric. In MOC the desired characteristics
543 of the simulation, which can represent either established biological data to be reproduced
544 or some desired hypothetical simulation outcome, are expressed as distributions. Import-
545 tantly, several such characteristics can be expressed, and MOC identifies those sets of
546 parameter values that deliver optimal tradeoffs against each.

547 We evaluated MOC on a well understood simulation, using it to reproduce a pre-
548 vious manual calibration effort and therein delivering a solution that was known to be
549 possible. The ARTIMMUS simulation was used, which simulates a mouse multiple scler-
550 osis disease model [18]. MOC delivered around 90 unique parameter value combinations,

551 each of which provided an optimal tradeoff in performance against the 5 target ARTIM-
552 MUS characteristics specified. This range of possible calibration solutions was unknown
553 *a priori*; the previous manual calibration of ARTIMMUS having delivered only one such
554 solution [11]. It would ordinarily fall on the simulation user to select one solution (set
555 of parameter values) to adopt as a baseline for subsequent simulation experimentation.
556 We investigated the significance of selecting solutions representing different extremes of
557 tradeoffs in delivering target simulation characteristics. A previous experiment with AR-
558 TIMMUS determined that removing the spleen, an important immune system organ,
559 resulted in exacerbated autoimmune symptoms. The results of re-performing this exper-
560 iment with different MOC solutions adopted as baseline parameter values were broadly,
561 but not absolutely, similar. Hence, adopting different calibration solutions can lead to
562 different experimental conclusions. It a strength of MOC that this range of solutions is
563 made explicit. Where possible, we recommend that MOC solutions be evaluated against
564 biological data to discard those that represent biologically unrealistic parameter values or
565 behaviours. Where this is not possible, we advocate performing experiments in replicate
566 using multiple MOC solutions such that the full range of possible results be established
567 before conclusions are drawn.

568 We demonstrated MOC in deriving simulation behaviours that were not known pos-
569 sible *a priori*: varying degrees of persisting autoimmunity in ARTIMMUS. MOC can be
570 applied to both parameters and initial conditions, at the same time, as demonstrated in
571 these calibration exercises. We do not consider simulation parameter values and initial
572 conditions as independent; a poor selection of initial condition values coupled with appro-
573 priate parameter values can still fail to deliver the desired simulation dynamic. MOC's
574 successful delivery of these previously unknown simulation dynamics presents an inter-
575 esting use case for MOC. It could be used to identify which parameters, and hence com-
576 ponents and pathways, need be manipulated to resolve a simulated disease state, therein
577 highlighting candidate therapeutic targets. Furthermore, for disease simulations that in-
578 corporate potential interventions, MOC can be used to determine optimal intervention
579 strategies that exploit synergies between several treatment options.

580 We surmise that MOC can support model selection and development. Accurately
581 simulating a biological system requires both an appropriate model of the biology, and ap-
582 propriate parameter values for that model. There typically exist several options for how

583 to represent a biological concept in simulation, the most suitable of which is often unclear.
584 Models must strike a balance between including sufficient complexity to accurately reflect
585 the biology's dynamics, whilst remaining sufficiently simplistic to offer insight. The un-
586 successful calibration of a given model of the biology can lead to two conclusions; first,
587 that the calibration process was simply unsuccessful in finding a solution that does exist,
588 a risk we argue is greatly lessened through MOC; or second, that the model is incapable
589 of replicating the biological dynamics in question. In this latter case, MOC can inform
590 simulation design, where a succession of putative models can be evaluated until calibra-
591 tion is successful. The possibility of directly applying MOC to the space of biological
592 abstractions, rather than parameter values, is intriguing, though extremely challenging
593 technically. Here, MOC would search for which cells were represented, and how. This
594 would encompass their interactions with one another, opting to ignore some found to
595 be irrelevant to the biological phenomenon of interest, or *vice versa*. The level of detail
596 through which molecular secretions and expressions were represented could also be de-
597 termined; is variable expression level necessary, or does simply 'present' vs 'not' suffice?
598 The challenge herein lies in building an agent-based simulation infrastructure capable
599 of capturing all these possibilities, and allowing the automated optimisation process to
600 manipulate them. The aforementioned point still applies, for each possible model, the
601 space of parameter values must also be investigated, as an accurate reflection of biology
602 requires both an appropriate model and corresponding parameter values. Hence, MOC
603 would be applied in a nested fashion, firstly over the space of biological representations,
604 and therein over the space of parameter values for each model.

605 Although our present investigation has employed an agent-based simulation, MOC
606 is applicable to other simulation paradigms also, such as ordinary differential equations
607 (ODE). Application to non-stochastic simulations, such as ODEs, requires significantly
608 less computational power, as there is no need to obtain simulation replicates in assessing
609 a candidate solution's fitness. We note that, from our experience in building them, not
610 all biological simulations are as computationally costly to execute and calibrate as AR-
611 TIMMUS. Each MOC calibration exercise has taken up to a week of time on a dedicated
612 computational facility. In this regard, terminating the MOC process when a threshold
613 level of overfitting is detected is pertinent (see Figure 8). Overfitting was detected in
614 all three of our ARTIMMUS parameter recalibration exercises, and selecting a threshold

615 of 0.5, wherein half of the MOC solutions at a given point no longer represent optimal
616 performance tradeoffs in an independent test, could as much as halve the computational
617 effort required.

618 The ability to detect overfitting in a multi-objective context is a novel contribution
619 of this work. Though a common strategy for stopping a single-objective optimisation
620 process, it was previously unclear how to deploy this strategy in a multi-objective context
621 [31]. There is another condition under which we feel it pertinent to terminate the MOC
622 process. The goal of MOC is to find parameter values yielding simulation dynamics that
623 closely resemble some target. As this alignment increases, and differences in solutions'
624 simulation performances reduce, it is possible that seemingly better alignments in fact
625 represent sampling artefacts arising from the stochastic simulation, rather than genuinely
626 superior parameter values. We note that detecting this in a statistically robust manner
627 is challenging, and as such we highlight it as potential further work.

628 This work fits within the context of a wider framework for supporting complex system
629 simulation, the CoSMoS framework [32]. CoSMoS advocates explicitly recording, typically
630 through graphical modelling [33], how real world concepts are translated into computer
631 code, and the implicit assumptions therein. In this context, MOC can help in relating
632 simulation results to biological data. The case where a distribution of results emerges
633 from a given biological experiment, even to the point where replicates or individuals
634 within an experiment exhibit completely different outcomes, can be handled in MOC
635 by defining bi-modal (or multi-modal) target distributions. A scenario wherein MOC
636 unexpectedly delivers several distinct and unconnected simulation phenotypes, rather
637 than a continuum of points on the Pareto front, is interesting. This can either suggest
638 the existence of additional phenotypes to look for in the biology, or if this can be ruled
639 out, suggests instead that the model being calibrated fails to accurately capture the
640 biology. This later case is an example of how MOC could drive simulation design and
641 development, as covered above. Related work on supporting the link of simulation to
642 biology proposes the construction of an argument wherein a claim such as 'this simulation
643 is an adequate representation of the biology' is supported by explicitly cited evidence [34].
644 In this context, application of MOC can raise confidence that appropriate parameter and
645 initial condition values have been identified. The range of possible values can be contrasted
646 against biological literature and data, excluding those deemed implausible. Subsequent

647 simulation experiments can be performed in replicate with those that remain, therein
648 highlighting the full range of results that are plausible in absence of better reason to rule
649 out particular parameter values. We argue that drawing conclusions from this nature of
650 simulation experimentation, and making explicit the full range of parameter values that
651 satisfy the calibration problem, leads to more robust conclusions.

652 In summary, our novel application of multi-objective optimisation in MOC presents
653 the multi-objective optimisation community with a new field of application, and one we
654 feel has considerable scope for growth. Importantly, it provides fundamental support for
655 a critical aspect of simulation-based biological experimentation: identifying parameter
656 values and initial conditions that align simulations with a complex target behaviour.

657 **Competing interests**

658 We have no competing interests.

659 **Authors' contributions**

660 MR conceived and designed the study, carried out all data generation and analysis re-
661 ported herein, and drafted the manuscript. KA conceived the study and drafted the
662 manuscript. LR conceived the study. JT conceived the study and drafted the manuscript.

663 **Acknowledgment**

664 We thank Dr. Fabian Held, Prof. Susan Stepney and Dr. Simon Poulding for engaging
665 discussions on statistics.

666 **Funding**

667 MR is supported by the David & Judith Coffey Life Lab. JT is part-funded by The Royal
668 Society.

References

- [1] Bauer AL, Beauchemin CA, Perelson AS. Agent-Based Modeling of Host-Pathogen Systems: The Successes and Challenges. *Inf Sci (Ny)*. 2009;179(10):1379–1389.
- [2] Harris TH, Banigan EJ, Christian Da, Konradt C, Tait Wojno ED, Norose K, et al. Generalized Lévy Walks and the Role of Chemokines in Migration of Effector CD8+ T Cells. *Nature*. 2012;486(7404):545–548.
- [3] An G, Mi Q, Dutta-Moscato J, Vodovotz Y. Agent-Based Models in Translational Systems Biology. *Wiley Interdiscip Rev Syst Biol Med*. 2009;1(2):159171.
- [4] Pienaar E, Dartois V, Linderman JJ, Kirschner DE. In Silico Evaluation and Exploration of Antibiotic Tuberculosis Treatment Regimens. *BMC Systems Biology*. 2015;9(1):79.
- [5] Thorne BC, Bailey AM, DeSimone DW, Peirce SM. Agent-Based Modeling of Multicell Morphogenic Processes during Development. *Birth Defects Research Part C, Embryo Today: Reviews*. 2007;81(4):344–353.
- [6] Senior AM, Charleston MA, Lihoreau M, Buhl J, Raubenheimer D, Simpson SJ. Evolving Nutritional Strategies in the Presence of Competition: a Geometric Agent-Based Model. *PLoS Computational Biology*. 2015;11(3):e1004111.
- [7] Kitano H. Systems Biology: a Brief Overview. *Science*. 2002;295(5560):1662–4.
- [8] Kitano H. Computational Systems Biology. *Nature*. 2002;420:206–210.
- [9] Cosgrove J, Butler J, Alden K, Read M, Kumar V, Cucurull-Sanchz L, et al. Agent-Based Modelling in Systems Pharmacology. *CPT: Pharmacometrics & Systems Pharmacology*. 2015;4(11):615629.
- [10] Dong C, Martinez GJ. T cells: the Usual Subsets; 2010. *Nature Reviews Immunology* poster: www.nature.com/nri/posters/tcellsubsets.
- [11] Read M, Andrews PS, Timmis J, Kumar V. Techniques for Grounding Agent-Based Simulations in the Real Domain: a case study in Experimental Autoimmune Encephalomyelitis. *Mathematical and Computer Modelling of Dynamical Systems*. 2012;18(1):67–86.

- 697 [12] Calvez B, Hutzler G. Automatic Tuning of Agent-Based Models Using Genetic Al-
698 gorithms. In: Sichman JS, Antunes L, editors. Multi-Agent-Based Simulation VI.
699 Springer Berlin Heidelberg; 2006. p. 41–57.
- 700 [13] Fabretti A. On the Problem of Calibrating an Agent Based Model for Financial
701 Markets. *Journal of Economic Interaction and Coordination*. 2012;8(2):277–293.
- 702 [14] Burnham KP, Anderson DR. Model Selection and Multimodel Inference: A Practical
703 Information-Theoretic Approach (2nd ed). vol. 172; 2002.
- 704 [15] Cohen IR. *Tending Adam’s Garden : Evolving the Cognitive Immune Self*. Elsevier
705 Academic Press; 2004.
- 706 [16] Kitano H. Biological Robustness. *Nature Reviews Genetics*. 2004;5(11):826–837.
- 707 [17] Deb K, Pratap A, Agarwal S, Meyarivan T. A Fast and Elitist Multiobjective
708 Genetic Algorithm: NSGA-II. *IEEE Transactions on Evolutionary Computation*.
709 2002;6(2):182–197.
- 710 [18] Read M, Andrews PS, Timmis J, Williams RA, Greaves RB, Sheng H, et al. Deter-
711 mining Disease Intervention Strategies using Spatially Resolved Simulations. *PLOS*
712 *ONE*. 2013;8(11):e80506.
- 713 [19] Baxter AG. The Origin and Application of Experimental Autoimmune En-
714 cephalomyelitis. *Nature Reviews Immunology*. 2007;7(11):904–912.
- 715 [20] Pachner AR. Experimental Models of Multiple Sclerosis. *Current Opinion in Neu-*
716 *rology*. 2011;24(3):291–299.
- 717 [21] Kumar V, Sercarz E. An Integrative Model of Regulation Centered on Recognition
718 of TCR Peptide/MHC Complexes. *Immunol Rev*. 2001;182:113–121.
- 719 [22] Kumar V. Homeostatic Control of Immunity by TCR Peptide-Specific Tregs. *J Clin*
720 *Invest*. 2004;114(9):1222–1226.
- 721 [23] Greaves RB, Read M, Timmis J, Andrews PS, Butler JA, Gerckens BO, et al. In
722 Silico Investigation of Novel Biological Pathways: the Role of CD200 in Regulation
723 of T Cell Priming in Experimental Autoimmune Encephalomyelitis. *Biosystems*.
724 2013;112(2):107–121.

- 725 [24] Williams RA, Greaves R, Read M, Timmis J, Andrews PS, Kumar V. In Silico
726 Investigation into Dendritic Cell Regulation of CD8Treg Mediated Killing of Th1
727 Cells in Murine Experimental Autoimmune Encephalomyelitis. *BMC Bioinformatics*.
728 2013;14 Suppl 6:S9.
- 729 [25] Vargha A, Delaney HD. A Critique and Improvement of the CL Common Language
730 Effect Size Statistics of McGraw and Wong. *Journal of Educational and Behavioral*
731 *Statistics*. 2000;25(2):101–132.
- 732 [26] Ben-Nun A, Ron Y, Cohen IR. Spontaneous Remission of Autoimmune En-
733 cephalomyelitis is Inhibited by Splenectomy, Thymectomy or Ageing. *Nature*.
734 1980;288(5789):389–390.
- 735 [27] Levine S, Sowinski R. Experimental Allergic Encephalomyelitis in Inbred and Out-
736 bred Mice. *Journal of Immunology*. 1973;110(1):139–43.
- 737 [28] Gold R. Understanding Pathogenesis and Therapy of Multiple Sclerosis via An-
738 imal Models: 70 Years of Merits and Culprits in Experimental Autoimmune En-
739 cephalomyelitis Research. *Brain*. 2006;129(8):1953–1971.
- 740 [29] Beeston T, Smith TR, Maricic I, Tang X, Kumar V. Involvement of IFN-g and
741 Perforin, but not Fas/FasL Interactions in Regulatory T Cell-Mediated Suppres-
742 sion of Experimental Autoimmune Encephalomyelitis. *Journal of Neuroimmunology*.
743 2010;15(229):91–97.
- 744 [30] Kumar V, Stellrecht K, Sercarz E. Inactivation of T Cell Receptor Peptide-
745 Specific CD4 Regulatory T Cells Induces Chronic Experimental Autoimmune En-
746 cephalomyelitis (EAE). *J Exp Med*. 1996;184(5):1609–1617.
- 747 [31] Dos Santos EM, Sabourin R, Maupin P. Overfitting Cautious Selection of Classifier
748 Ensembles with Genetic Algorithms. *Information Fusion*. 2009;10(2):150–162.
- 749 [32] Bown J, Andrews PS, Deeni Y, Goltsov A, Idowu M, Polack FAC, et al. Engineering
750 Simulations for Cancer Systems Biology. *Current Drug Targets*. 2012;13(12):1560–
751 1574.

752 [33] Read M, Andrews PS, Timmis J, Kumar V. Modelling Biological Behaviours with the
753 Unified Modelling Language: an Immunological Case Study and Critique. Journal
754 of the Royal Society Interface. 2014;11(9):20140704.

755 [34] Alden K, Andrews PS, Polack FAC, Veiga-fernandes H, Coles MC, Timmis J, et al.
756 Using Argument Notation to Engineer Biological Simulations with Increased Confi-
757 dence. Journal of the Royal Society Interface. 2015;12:20141059.

758 **Figure and table captions**

759 **Figure 1. The ARTIMMUS simulation, used as a testcase for evaluating MOC.**
760 **A**, the the major cell types represented in ARTIMMUS, and their key influences on one
761 another. Red and green arrows respectively indicate activities that perpetuate autoimmu-
762 nity or mediate recovery. Figure adapted from [11]. **B**, the baseline dynamic of ARTIM-
763 MUS, depicting four T cell population sizes over time. The simulation behaviour depicted
764 here forms a calibration target for MOC in Section 4. Lines correspond to like-coloured
765 cells in Figure A; these colours are maintained throughout the manuscript. Error bars
766 capture 90% of the data derived from 500 simulation executions, timeseries lines indicate
767 median population sizes at each time point.

768

769

770 **Figure 2. Overview of the Multi-Objective Calibration (MOC) concept.** **A**,
771 The desired (target) simulation dynamics are defined as distributions (only 2 shown): the
772 desired distributions of peak cell number and the times at which these occur. Distribu-
773 tions are depicted as histograms, or the corresponding cumulative distribution functions
774 describing the proportion of samples in the distribution (y axis) that hold a given value
775 or less (x axis). **B**, the capacity for putative simulation parameter (only 2 shown) values,
776 termed *candidate solutions*, to reproduce target dynamics is evaluated. The Kolmogorov-
777 Smirnov (KS) statistic quantifies the difference between target and a given candidate
778 solution's simulation performance (left); this metric is termed an *objective*. By sampling
779 and evaluating regions of parameter space we identify those that provide good alignment
780 with a given objective, illustrated through greyscale heatmaps (right). No single region
781 of parameter space maximizes performance against all objectives (only 2 shown), there

782 exist inherent tradeoffs. A heuristic (guided) search strategy, NSGA-II, is employed to
783 strategically sample parameter space. **C**, *solutions* representing optimal tradeoffs in per-
784 formance against each objective are identified, collectively termed the *Pareto front* (left).
785 These solutions are *Pareto-equivalent* (pink): no solution has been found that represents
786 an improvement in one objective without a worsening in another. Sub-optimal candidate
787 solutions are discarded (blue). Pareto-equivalent solutions may reside in disparate regions
788 of parameter space(right).

789

790

791 **Figure 3. Multi-objective calibration (MOC) successfully re-calibrates AR-**
792 **TIMMUS parameters against 5 objectives.** The best solution's, that with lowest
793 Λ score, target simulation dynamics are shown. The solution dataset comprises 200 sim-
794 ulation replicates, the target comprises 500. **A**, T cell population sizes over time, for
795 both target (dotted line) and solution (solid line). The median values from each dataset
796 at the given point in the time series are plotted. **B-F**, cumulative distribution functions
797 showing alignment of solution and target distributions of values for each objective, with
798 titles giving KS values. These graphs show the distribution of calibration target values
799 obtained in each dataset: the y-axis indicates the proportion of items in the distribution
800 holding a value less than or equal to the corresponding x-axis value. Objectives are: B,
801 peak CD4Th1 population size cell; C, time at which this peak occurs; D, CD4Th2 pop-
802 ulation size at 30 days; E, peak CD4Treg population size; F, peak CD8Treg population
803 size. These data represent the first of three independent recalibration experiments.

804

805

806 **Figure 4. Automated re-calibration of ARTIMMUS parameters delivers so-**
807 **lutions approximating the original manually-calibrated parameter values.** Box
808 plots are shown for each of three independent calibration exercises. The horizontal green
809 line represents the manually-calibrated parameter values. Calibration was performed over
810 5 objectives: the peak population sizes of Th1, CD4Treg, CD8Treg cells, the time at which
811 the Th1 population peaks, and the number of Th2 cells at 30 days. Parameters subject
812 to calibration are listed in Table 1, see Figure S1 for an explanation of their operation in
813 ARTIMMUS. Values shown above each plot are the Kolmogorov-Smirnov scores between

814 distributions, shown to one significant figure; the associated p-values are: *, $p < 0.01$ and
815 **, $p < 0.001$. Outliers in boxplots are defined as lying beyond the first or third quartiles
816 by 1.5 times the interquartile range.

817

818

819 **Figure 5. The range of objective values that constitute the Pareto front de-**
820 **ri-ved through MOC re-calibration of ARTIMMUS.** Box plots are shown for each of
821 three independent calibration exercises. These objective values correspond to the Pareto
822 front and associated ARTIMMUS parameter values of Figure 4. Calibration was per-
823 formed against five objectives: **A**, the peak population size of Th1 cells; **B**, the time at
824 which this occurred; **C**, the number of Th2 cells at 30 days; **D**, the peak population size
825 of CD4Treg cells; **E**, the peak population size of CD8Treg cells. Statistical and boxplot
826 formatting are as in Figure 4.

827

828

829 **Figure 6. Do different regions of MOC's Pareto front of solutions give rise**
830 **to different results in subsequent experimentation?** Top, an overview of the ex-
831 perimental procedure. **1**, Pareto front members representing objective value extremes
832 are identified (only two objectives shown in example). **2**, The simulation parameters
833 represented by such members are adopted in performing a control and splenectomy ex-
834 periment, with 200 replicate simulations in each group. **3**, Key performance indicators are
835 extracted from the resultant distributions of simulation dynamics, indicated here is the
836 peak CD4Treg population size within each individual simulation. **4**, Performance indica-
837 tors are statistically contrasted for splenectomy and control experiments. These statistics
838 are examined across different Pareto front members, thereby gauging the extent to which
839 experimental results critically depend on which Pareto-equivalent parameter values are
840 adopted in simulation. **Tables**, columns represent extreme Pareto front solutions, defined
841 as having either the highest or lowest KS value for each of the five objectives used in cal-
842 ibration (see Section 4). The objective KS value scores are shown in parentheses. Only
843 the first occurrence of each solution is shown, with subsequent entries indicated by '-'.
844 Rows indicate the difference between control and splenectomy simulations based on each
845 solution according to key indicators of simulation behaviour, as measured by the Vargha-

846 Delaney *A test* [25]. The original *A test* scores for the manually-calibrated simulation
847 are shown ('orig'), as is the biggest difference in *A test* score observed between manually-
848 and automatically-calibrated simulations ('diff'). Values highlighted in red represent four
849 differences in candidate and original *A test* scores that are notably larger than differences
850 observed elsewhere. 'D.c.' indicates 'direction change', where there exists at least one
851 candidate with for which the *A test* score lay on the other side of 0.5 from the original,
852 indicating that the distribution of values under splenectomy increased in the original ex-
853 periment but decreased for the candidate (or *vice versa*).

854

855

856 **Figure 7. Employing MOC to discover parameter and initial condition val-**
857 **ues delivering simulation dynamics not known to exist *a priori*: persisting**
858 **autoimmune states of varying severity.** Three severities are explored, represented
859 as columns. They are, a mean of 50, 200 or 500 Th1 cells at 60 days (with standard
860 deviations of 10, 100 and 200 respectively). A second objective is employed in all cases,
861 1000 Th1 cells at 15 days, which drives the establishment of autoimmunity. Each severity
862 is calibrated in three independent experiments, and shown here are the solutions exhibit-
863 ing lowest Λ values from a representative calibration of each experiment. The first row
864 of graphs depicts the median T cell time-series. The second row shows the candidate's
865 performance against an objective of 1000 Th1 at 60 days (standard deviation = 200). The
866 last row depicts the second objective, the (respective) number of T cells at 60 days.

867

868

869 **Figure 8. Terminating MOC when overfitting occurs.** Overfitting describes the
870 case when solutions generated by an optimisation process, e.g. MOC, better resemble
871 the target data than the underlying distribution from which it was drawn. **A**, in many
872 contexts, such as animal experiments, only limited samples of a phenomenon can be ob-
873 tained. The samples will broadly, but not exactly, reflect the underlying distribution. **B**,
874 an overfitted candidate solution more closely resembles the target data than the underly-
875 ing distribution from which the target data was drawn. Detecting this is difficult because
876 the true underlying distribution cannot be absolutely known. **C**, a common strategy in
877 single-objective optimisation problems is to divide the available data into two, a training

878 dataset and a validation dataset. The training dataset is used as the target in obtaining
879 successively better quality solutions. The validation dataset is used as an independent
880 check. Overfitting is detected when solutions more closely resemble the training dataset
881 than validation dataset. This is illustrated in **D**, where early solutions generally offer
882 improved performance against both datasets. It is only in later stages that solutions so
883 closely reflect the target dataset that they diverge from the validation dataset. This is
884 when the process should be stopped. **E**. Overfitting can be detected in multi-objective
885 optimisation, such as MOC, by maintaining Pareto fronts of optimal solutions against
886 both training and validation data independently. The degree of overfitting is reflected
887 in the proportion of training data Pareto front solutions that are not members of the
888 validation data Pareto front.

889

890

891 **Figure 9. Empirical results for detecting overfitting in MOC, and when to**
892 **terminate the process accordingly.** We generated a validation dataset using AR-
893 TIMMUS’s previous manually-calibrated parameter values, and retrospectively analysed
894 how overfitted MOC solutions would have been on the three MOC calibration exercises
895 reported in Section 4. **A**, The overfittedness, defined as the proportion of MOC Pareto
896 front solutions that are not also members of a similar Pareto front maintained for the
897 validation data, at each MOC generation. **B**, the generation at which MOC would have
898 been terminated for a given overfittedness threshold value.

899

900

901 **Table 1.** The ARTIMMUS parameters (top) and initial conditions (bottom) subject to
902 calibration, their baseline (manually-calibrated) values, and the lower and upper bounds
903 of values they may be assigned during MOC.

904

905

906 **Table 2.** The best solution, being that with the lowest Λ value, arising from each of
907 three independent calibration exercises. Shown are each of the five objective KS values.
908 We independently investigated the calibration of both ARTIMMUS parameters (top) and
909 initial conditions (bottom). High quality calibrations, as indicated by low KS values, were

910 obtained in all cases.

911 **Figures and tables**

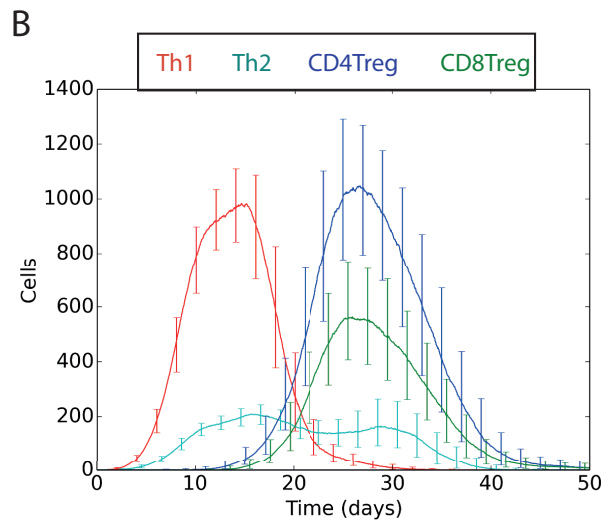
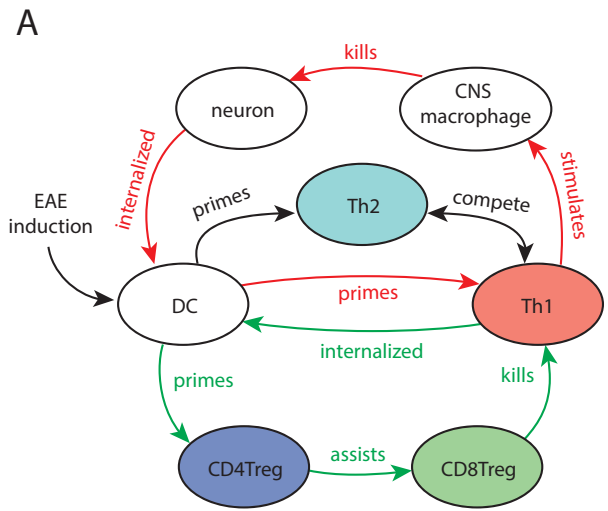
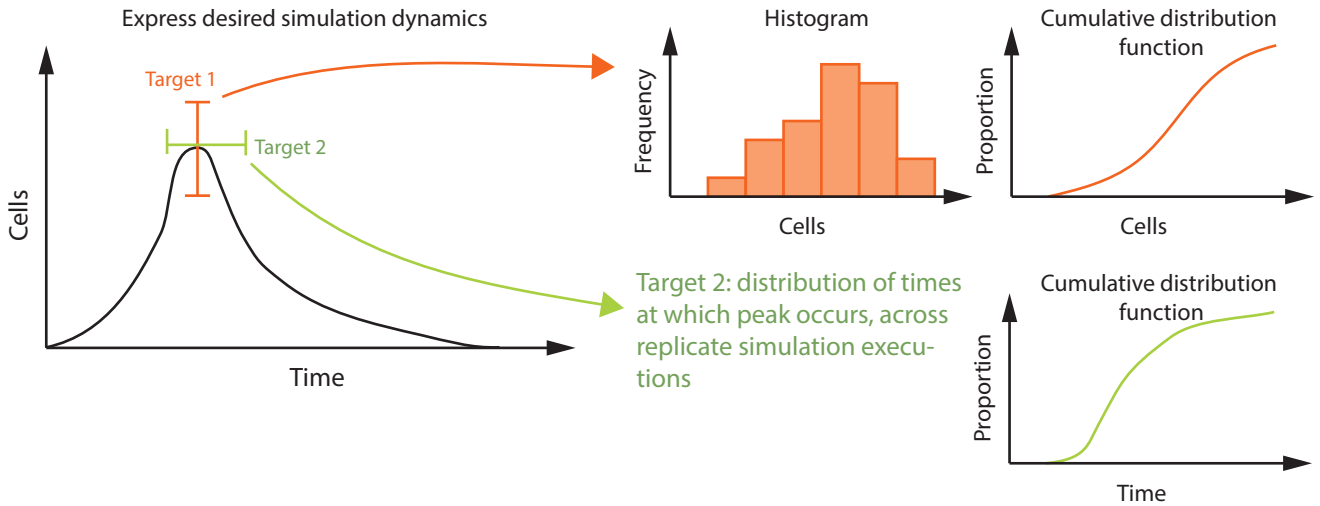
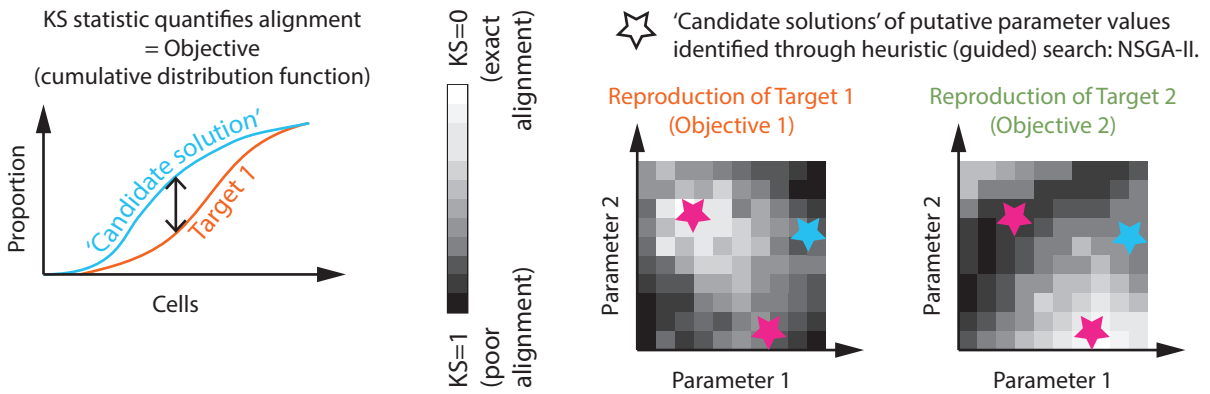


Figure 1:

A. Define target simulation dynamics



B. Evaluate putative simulation parameters' reproduction of target dynamic



C. Identify Pareto-optimal solutions (the Pareto front)



Figure 2:

Calibrate ARTIMMUS parameters, exercise 1

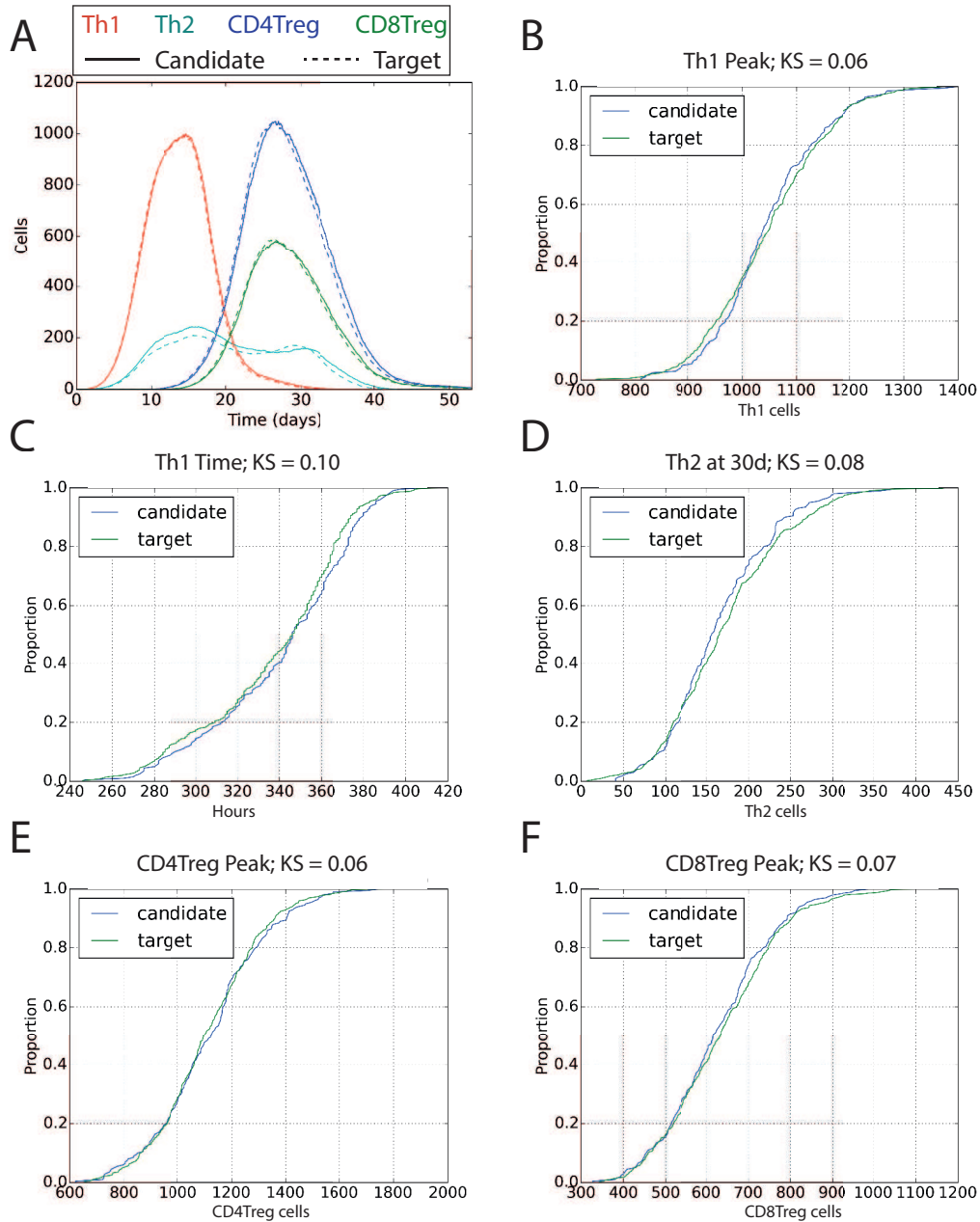


Figure 3:

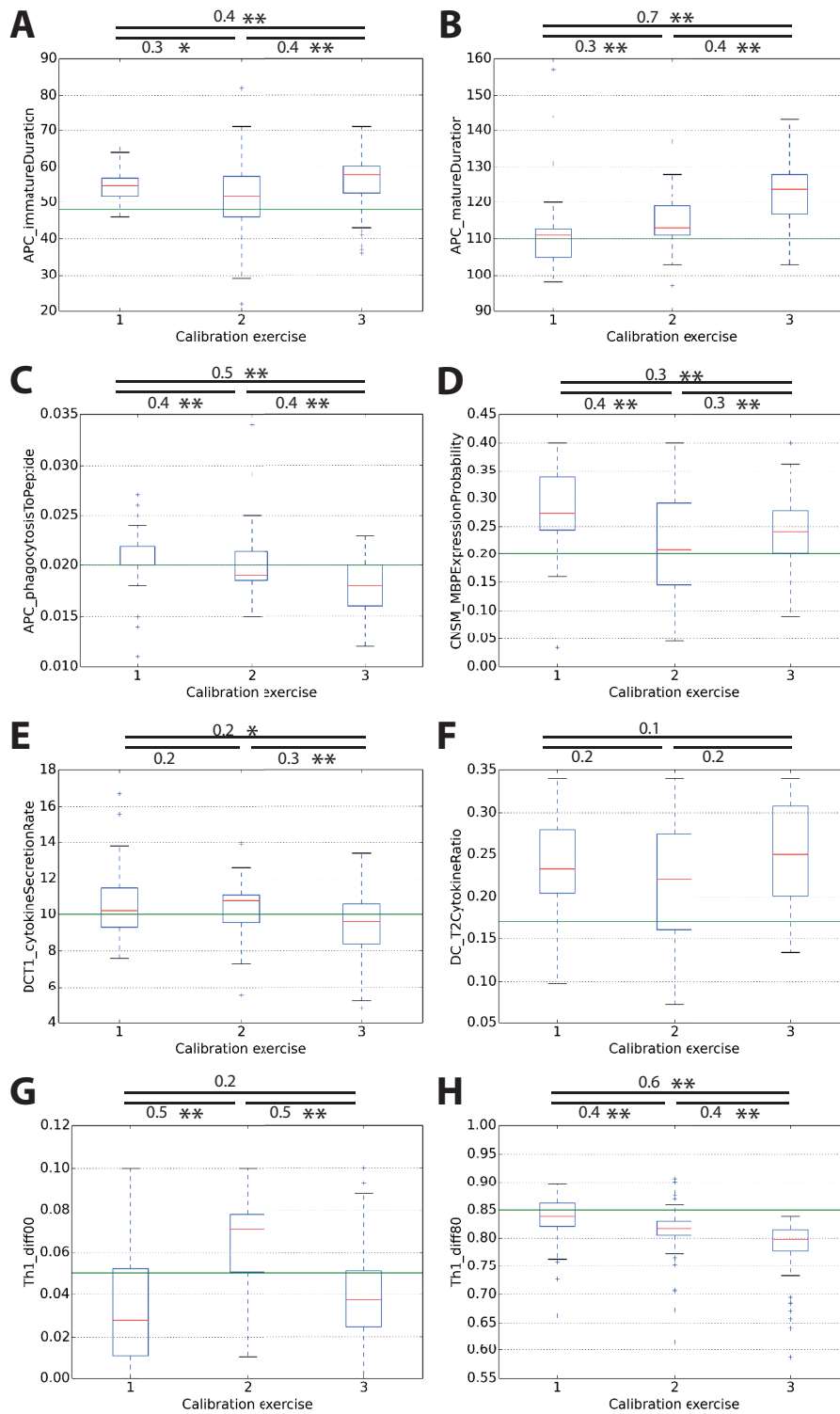


Figure 4:

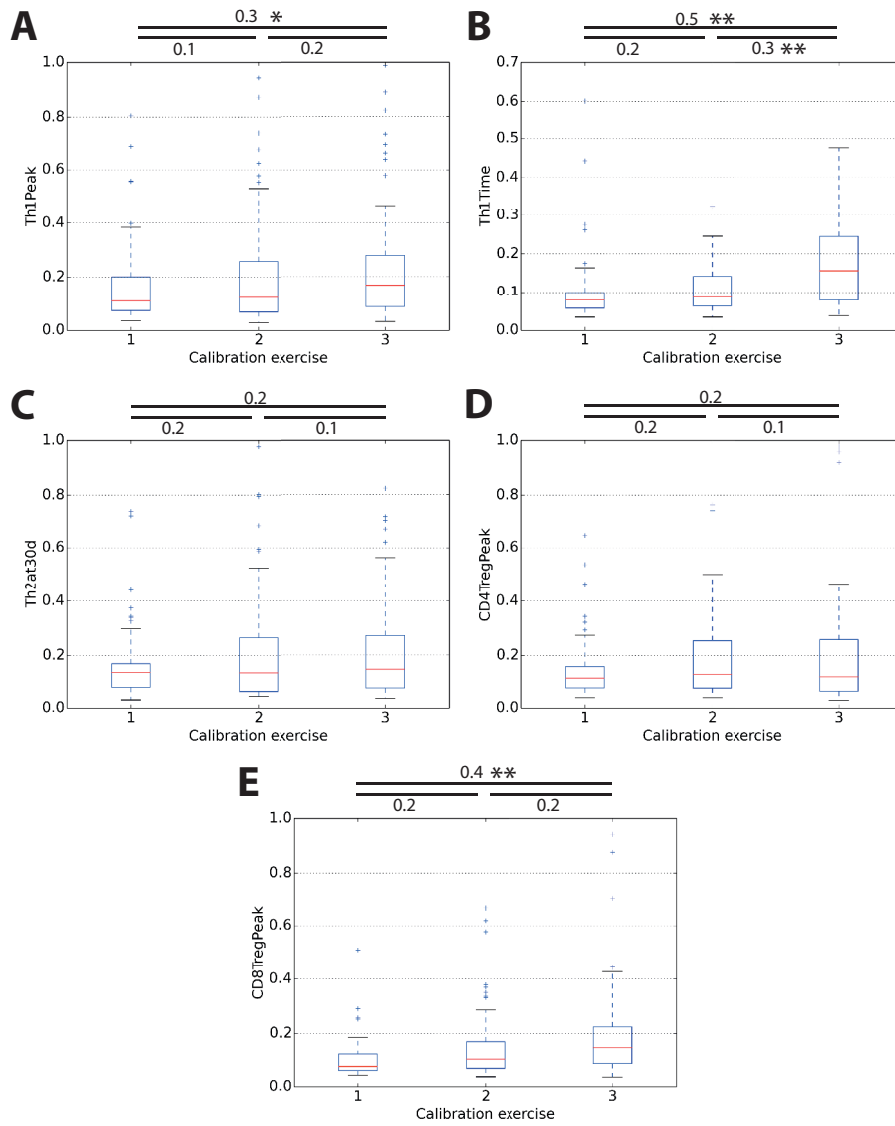
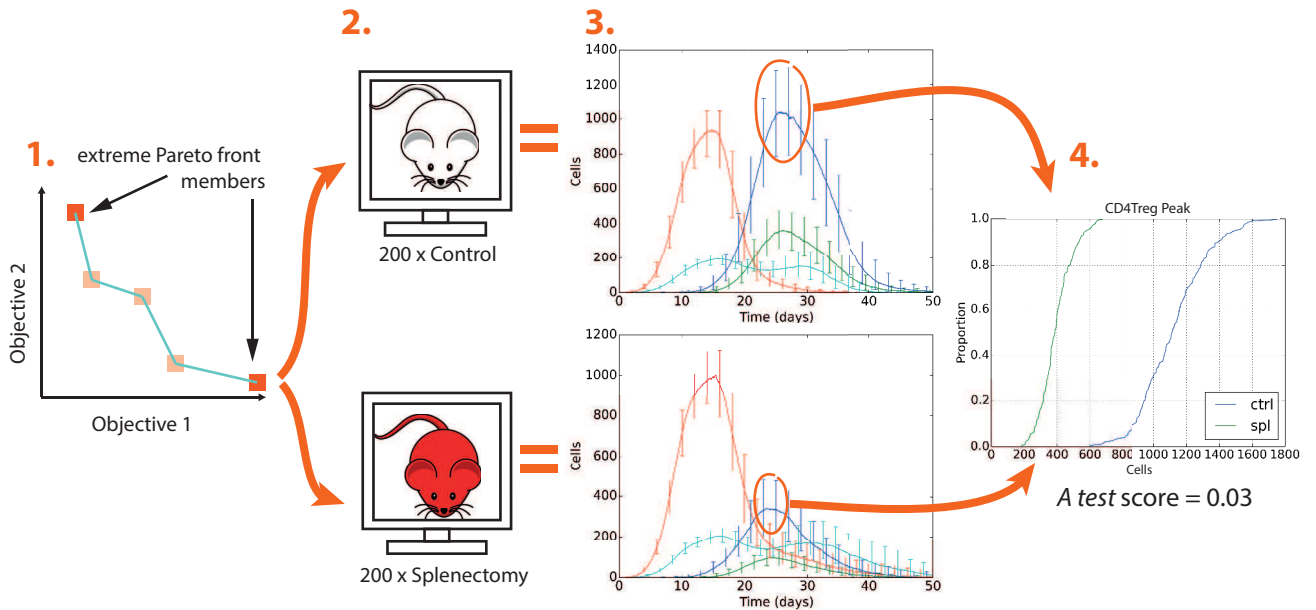


Figure 5:



Calibration Exercise 1, Vargha-Delaney A Test Scores													
Response	Th1Peak		Th1Time		Th2at30d		CD4TregPeak		CD8TregPeak		orig	diff	d.c.
	g31c32 (KS=0.04)	g11c33 (0.81)	g29c37 (0.03)	g10c13 (0.60)	g31c6 (0.03)	g23c60 (0.73)	g30c0 (0.04)	g23c22 (0.65)	g10c13 (0.04)	g23c22 (0.51)			
Th1 Peak	0.70	0.73	0.65	0.67	0.67	0.66	0.71	0.67	-	-	0.62	0.11	
Th1 Time	0.53	0.51	0.55	0.42	0.53	0.52	0.52	0.52	-	-	0.47	0.08	Y
Th2 Peak	0.60	0.74	0.66	0.66	0.67	0.67	0.61	0.62	-	-	0.66	0.08	
Th2 Time	0.47	0.53	0.68	0.46	0.51	0.56	0.53	0.47	-	-	0.58	0.11	Y
CD4Treg Peak	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	-	-	0.00	0	
CD4Treg Time	0.23	0.24	0.26	0.26	0.20	0.28	0.20	0.24	-	-	0.21	0.07	
CD8Treg Peak	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	-	-	0.00	0	
CD8Treg Time	0.24	0.24	0.29	0.27	0.23	0.25	0.25	0.26	-	-	0.23	0.06	
Th1 at 40d	0.98	0.96	0.92	0.89	0.96	0.58	0.96	0.94	-	-	0.95	0.37	

Calibration Exercise 2, Vargha-Delaney A Test Scores													
Response	g30c25 (KS=0.03)	g6c35 (0.94)	g30c34 (0.03)	g9c63 (0.33)	g15c14 (0.04)	g30c58 (0.98)	g17c61 (0.04)	g14c54 (0.77)	g9c56 (0.04)	g14c54 (0.67)	orig	diff	d.c.
	Th1 Peak	0.63	0.69	0.66	0.66	0.71	0.65	0.65	0.66	0.68			
Th1 Time	0.53	0.49	0.55	0.47	0.50	0.53	0.54	0.51	0.48	-	0.47	0.08	Y
Th2 Peak	0.66	0.67	0.65	0.70	0.68	0.39	0.65	0.61	0.70	-	0.66	0.27	Y
Th2 Time	0.56	0.51	0.63	0.63	0.50	0.31	0.54	0.48	0.54	-	0.58	0.27	Y
CD4Treg Peak	0.00	0.00	0.00	0.01	0.00	0.00	0.00	0.00	0.00	-	0.00	0.01	
CD4Treg Time	0.27	0.30	0.23	0.19	0.20	0.21	0.17	0.22	0.29	-	0.21	0.09	
CD8Treg Peak	0.00	0.00	0.00	0.01	0.00	0.00	0.00	0.00	0.00	-	0.01	0.01	
CD8Treg Time	0.27	0.26	0.23	0.17	0.21	0.20	0.22	0.23	0.28	-	0.23	0.06	
Th1 at 40d	0.96	0.56	0.94	0.91	0.94	1.00	0.97	0.97	0.82	-	0.95	0.39	

Calibration Exercise 3, Vargha-Delaney A Test Scores													
Response	g5c39 (KS=0.03)	g23c10 (0.99)	g28c47 (0.04)	g13c21 (0.48)	g11c6 (0.04)	g5c39 (0.82)	g24c62 (0.03)	g23c10 (0.98)	g23c40 (0.04)	g23c10 (0.94)	orig	diff	d.c.
	Th1 Peak	0.65	0.65	0.63	0.70	0.68	-	0.65	-	0.70			
Th1 Time	0.43	0.49	0.49	0.42	0.47	-	0.47	-	0.52	-	0.47	0.05	Y
Th2 Peak	0.67	0.64	0.70	0.68	0.67	-	0.67	-	0.67	-	0.66	0.04	
Th2 Time	0.59	0.51	0.52	0.50	0.47	-	0.62	-	0.57	-	0.58	0.11	Y
CD4Treg Peak	0.07	0.02	0.00	0.00	0.00	-	0.00	-	0.00	-	0.00	0.07	
CD4Treg Time	0.30	0.35	0.21	0.24	0.26	-	0.20	-	0.19	-	0.21	0.14	
CD8Treg Peak	0.04	0.01	0.00	0.00	0.00	-	0.00	-	0.00	-	0.00	0.04	
CD8Treg Time	0.27	0.40	0.17	0.26	0.27	-	0.20	-	0.21	-	0.23	0.17	
Th1 at 40d	0.93	0.83	0.97	0.90	0.92	-	0.97	-	0.95	-	0.95	0.12	

Figure 6:

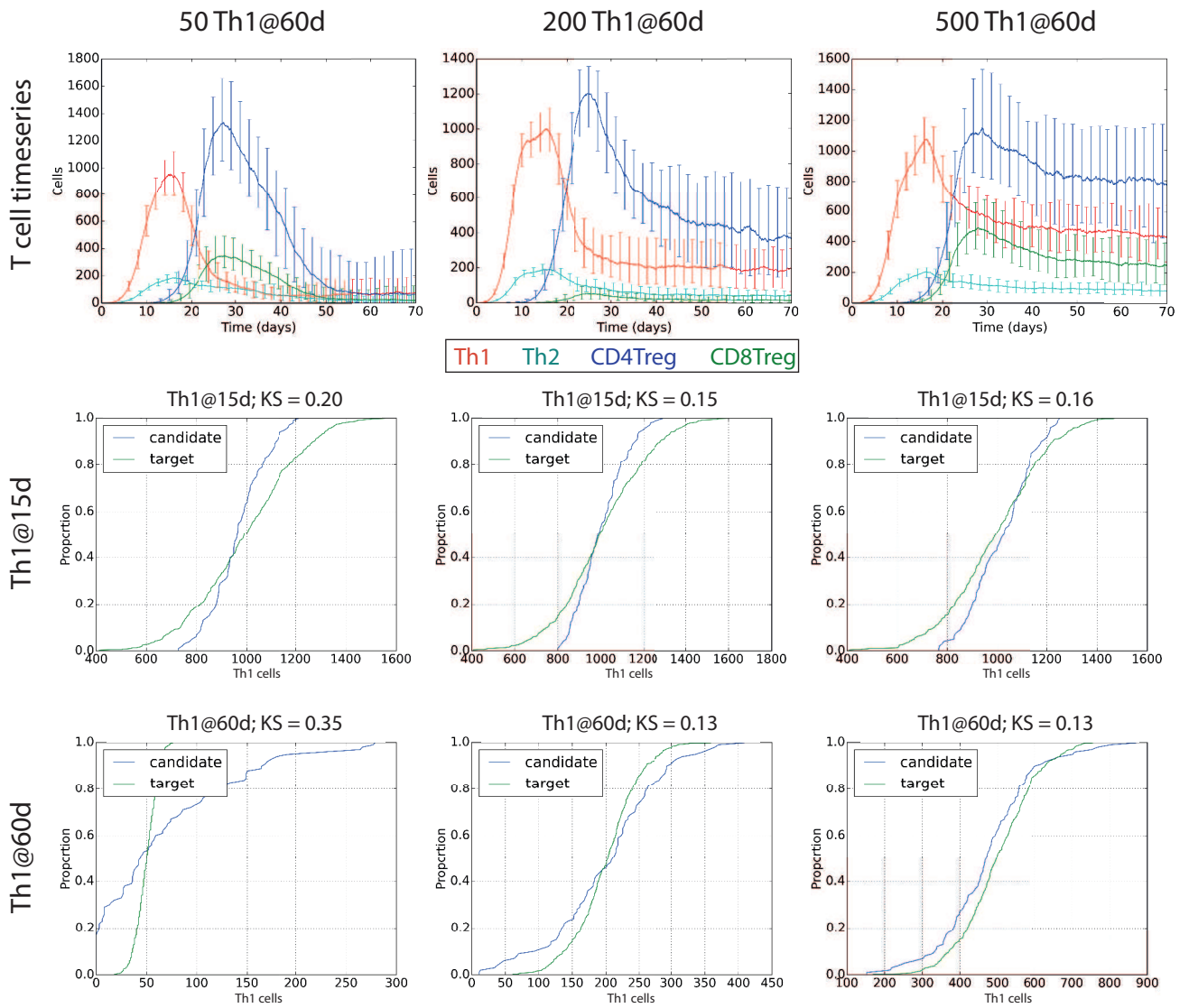


Figure 7:

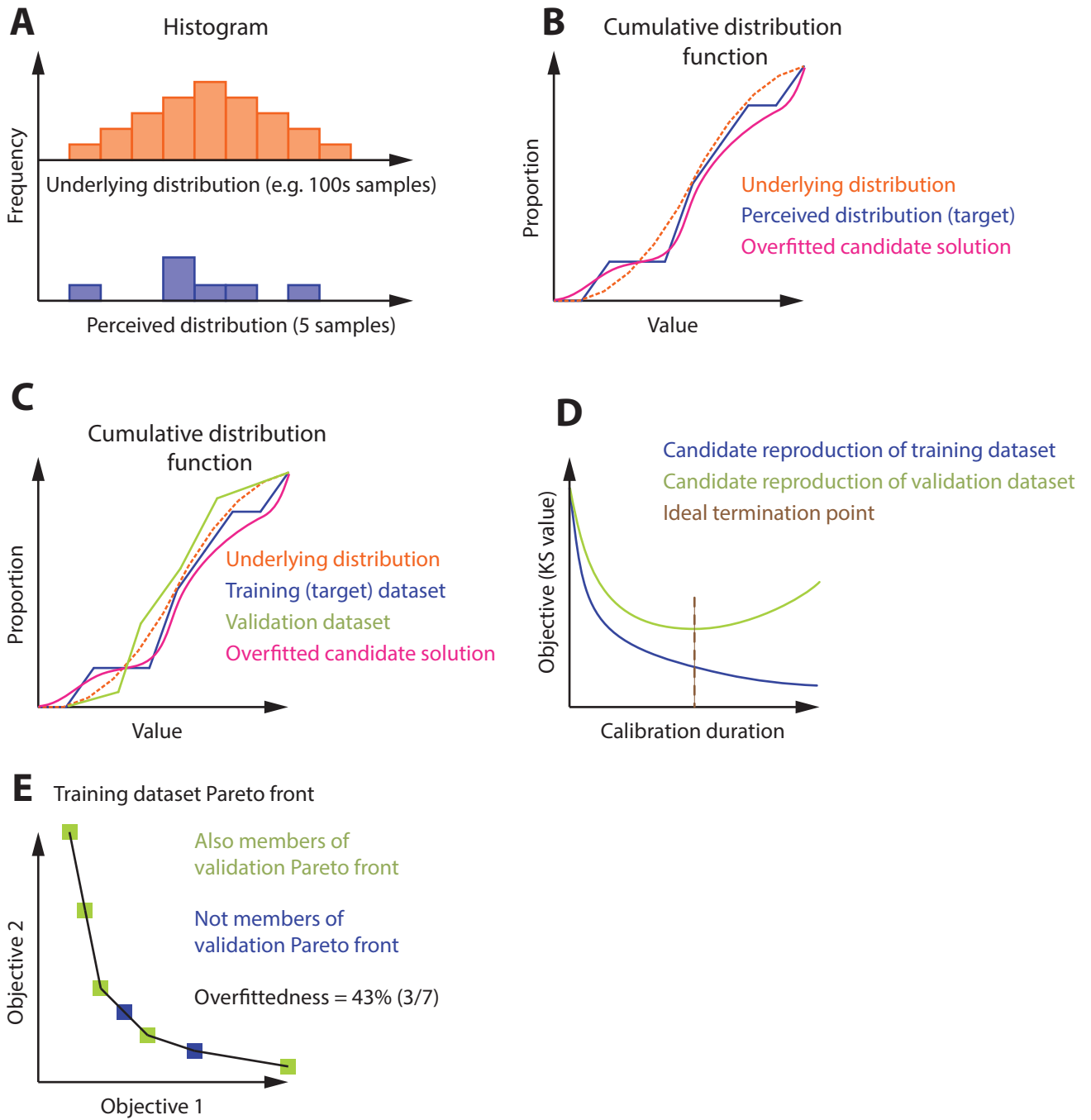


Figure 8:

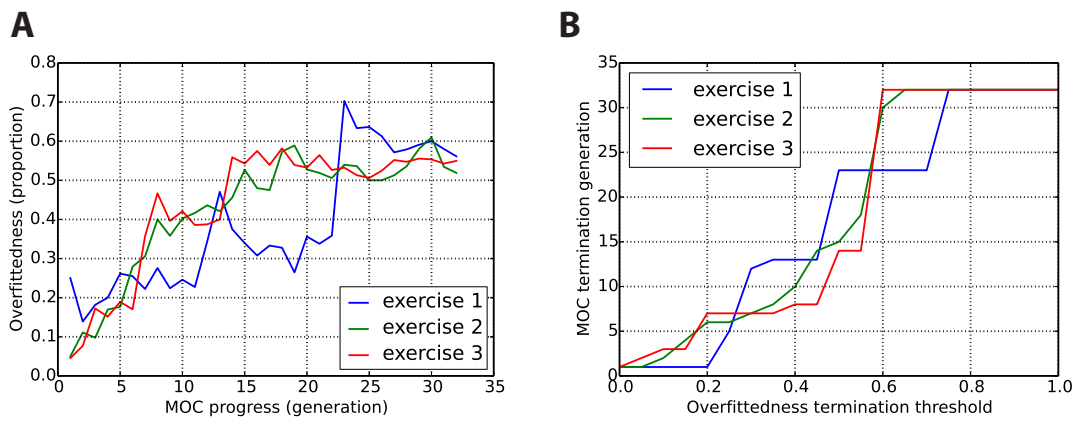


Figure 9:

Parameters calibrated			
Parameter	Baseline value	Lower bound	Upper bound
<i>APC_immatureDuration</i>	48	0	96
<i>APC_matureDuration</i>	110	0	220
<i>APC_phagocytosisToPeptide</i>	0.02	0	0.04
<i>CNSM_MBPEXpressionProbability</i>	0.2	0	0.4
<i>DCT1_cytokineSecretionRate</i>	10	0	20
<i>DC_T2CytokineRatio</i>	0.17	0	0.34
<i>Th1_diff00</i>	0.05	0	0.1
<i>Th1_diff80</i>	0.85	0	1.0
Initial conditions calibrated			
Initial condition	Baseline value	Lower bound	Upper bound
<i>numTh</i>	40	0	80
<i>numCD4Treg</i>	30	0	60
<i>numCD8Treg</i>	30	0	60
<i>numCNS</i>	500	0	1000
<i>numCNSMacrophage</i>	75	0	150
<i>numDC</i>	10	0	20
<i>numDCCNS</i>	40	0	80
<i>numDCSpleen</i>	100	0	200

Table 1:

Calibration on parameters					
Calibration exercise	Objective KS value				
	Th1Peak	Th1Time	Th2at30d	CD4TregPeak	CD8TregPeak
1	0.06	0.10	0.08	0.06	0.07
2	0.08	0.06	0.06	0.05	0.07
3	0.05	0.08	0.14	0.08	0.05
Calibration on initial conditions					
Calibration exercise	Objective KS value				
	Th1Peak	Th1Time	Th2at30d	CD4TregPeak	CD8TregPeak
1	0.06	0.08	0.04	0.03	0.06
2	0.04	0.08	0.10	0.11	0.12
3	0.06	0.06	0.06	0.07	0.05

Table 2: