

History matching with subset simulation

Z. T. Gong^a, F. A. DiazDelaO^b, P. O. Hristov^{c,*}, M. Beer^d

^a*CRRC Qingdao Sifang Co., LTD., China.*

^b*Clinical Operational Research Unit, Department of Mathematics, University College London, London WC1H 0BT, UK*

^c*Institute for Risk and Uncertainty, School of Engineering, University of Liverpool, Liverpool L69 7ZF, UK*

^d*Institute for Computer Science in Civil Engineering, Leibniz University, Hannover 30167, Germany*

Abstract

Computational cost often hinders the calibration of complex computer models. In this context, history matching is becoming a widespread calibration strategy, with applications in many disciplines. History matching (HM) uses a statistical approximation - also known as an emulator - to the model output, in order to mitigate computational cost. The process starts with an observation of a physical system. It then produces progressively more accurate emulators to determine a non-implausible domain: a subset of the input space that provides a good agreement between the model output and the data, conditional on the model structure, the sources of uncertainty and an implausibility measure. In HM, it is essential to generate samples from the non-implausible domain, in order to run the model and train the emulator until a stopping condition is met. However, this sampling can be very challenging, since the non-implausible domain can become orders of magnitude

*Corresponding author.

Email address: p.hristov2@liv.ac.uk (P. O. Hristov)

smaller than the original input space very quickly. This paper proposes a solution to this problem using subset simulation, a rare event sampling technique that works efficiently in high dimensions. The proposed approach is demonstrated via calibration and robust design examples from the field of aerospace engineering.

Keywords: History matching, Subset simulation, Gaussian process emulation, Robust design.

1 Introduction

The use of computer models (also known as simulators) to study complex systems and environments is indispensable in modern scientific research. The reliability of these models depends critically on how well they are calibrated to experimental data. Otherwise, model-based decisions run the risk of being misguided or ill-informed. One of the challenges of model calibration is that several sources of uncertainty must be taken into account. This uncertainty originates (for instance) due to process idealisations, model assumptions and computational cost. In order to provide evidence of predictive reliability, it is essential that any model is calibrated taking into account these sources of uncertainty.

Typically, high-fidelity computer models of complex phenomena are computationally expensive. In the context of uncertainty quantification, this characterisation usually describes models whose evaluation time prohibits their repeated use in any form of sampling-based analysis. This feature presents a challenge to classical calibration techniques, which require a considerable number of simulator runs to identify an acceptable match between

18 model and data. Furthermore, the analyst might face an added challenge
19 if not only the model, but also the generation of experimental data, is ex-
20 pensive or unfeasible. Despite the importance and necessity of efficient cali-
21 bration methods for complex computer codes, their development has lagged
22 behind their application [1]. Instead, simple goodness-of-fit measures such
23 as distance-based methods and least squares (see *e.g.* [2]) or likelihood func-
24 tions (consult [3] and references therein) are applied. Neither of these may
25 be suitable when computational cost and high dimensional input are consid-
26 ered. This is due to the fact that goodness-of-fit measures typically require
27 large data sets to achieve a reliable quantification of the degree of agreement
28 between observations and simulator realizations. Likewise, the likelihood of
29 complex simulators is usually intractable and approximations may be re-
30 quired (see *e.g.* [4]).

31 History matching (HM) [5, 6] is a form of calibration for complex and com-
32 putationally expensive numerical models. It uses Bayesian emulation [7] to
33 tackle computational cost. Emulation means building a statistical approx-
34 imation to the original simulator, thus allowing affordable inference about
35 its output. History matching also defines an implausibility measure, which
36 is used to reduce the input space by finding an input subspace that pro-
37 vides a reasonable match between the model output and experimental data,
38 given the model structure and various sources of uncertainty. This input
39 space reduction is achieved by building progressively more accurate emula-
40 tors, which in practice results in HM becoming an iterative process. The
41 resulting input subspace is known in the literature as non-implausible do-
42 main, non-implausible set or Not-Ruled-Out-Yet (NROY) space [8]. History

43 matching has been successfully applied in epidemiology [1], galaxy formation
44 modeling [8], oil reservoir analysis [9] and large climate systems modelling
45 [10], amongst many other applications.

46 The sequential generation of samples from the non-implausible domain
47 at every HM iteration has remained an open and complex problem. This
48 mainly stems from the fact that the non-implausible domain can be orders
49 of magnitude smaller than the original input space [11]. A notable example
50 of a field of study in which a similar challenge is encountered is *engineering*
51 *reliability analysis*. The main aim of this type of reliability analysis is to
52 identify the conditions under which a physical system fails. In that context,
53 failure means that the demand has exceeded the capacity of the system, ac-
54 cording to a model of the system and a criterion guided by expert knowledge.
55 Reliability analysis aims at generating samples from the *failure set*, that is,
56 the set of model input configurations that lead to failure. This allows the
57 characterisation of different modes in which the system can fail and to esti-
58 mate the probability of failure. If an engineering system is well-designed and
59 the model is a good representation of the system, the volume of the failure
60 domain is expected to be orders of magnitude smaller compared to the input
61 space. Since this can also be the case for the non-implausible domain within
62 HM, this opens the prospect of treating it as if it were a failure set within
63 reliability analysis.

64 Subset simulation (SuS) [12] is a widely used technique in engineering re-
65 liability computations and rare event simulation. Unlike direct Monte Carlo,
66 SuS models a rare event, which has a small failure probability, as contained
67 in a nested sequence of less-rare events. Eventually, the probability of failure

68 can be computed as the product of larger conditional probabilities given the
69 occurrence of each preceding event. Markov chain Monte Carlo (MCMC)
70 is used to generate the conditional samples that belong to the intermediate
71 failure events. Based on this strategy, SuS generates samples selectively, to
72 efficiently populate the target failure set. Given the potential disparity in
73 the size of the original input space and the non-improbable domain (in the
74 context of HM); and the potential disparity in the size of the original input
75 space and the failure domain (in the context of reliability analysis), this pa-
76 per proposes the use of SuS as an efficient sampler of the non-improbable
77 domain within each wave of HM.

78 The remainder of the paper is organized as follows. Section 2 presents an
79 overview of HM. Section 3 reviews the details of SuS. Section 4 presents the
80 proposed approach, in which SuS is used to sample from the non-improbable
81 domain in HM. The resulting procedure is demonstrated in a calibration
82 context in Section 5 and in an industrial context for robust aircraft design in
83 Section 6. Finally, section 7 provides some conclusions.

84 **2. History matching**

85 *2.1. Procedure overview*

86 A rigorous description of the relationship between a model and the un-
87 derlying physical process requires the identification and inclusion of different
88 sources of uncertainty. Let y denote the true value of the physical process. A
89 modeller analysing the process can only observe a noisy version of this value.
90 Let $z = y + \epsilon_{me}$ be this noisy observation, where ϵ_{me} is measurement error.
91 This type of error, also called *observational uncertainty*, can be thought of

92 as a random variable with zero mean and finite variance. The modeller then
 93 represents the physical process through a numerical simulator, which defines
 94 an input-output mapping $\eta : \mathbb{R}^d \rightarrow \mathbb{R}$. Let $\eta(\mathbf{x})$ denote the simulator output
 95 as a function of some input vector $\mathbf{x} \in \mathcal{X} \subseteq \mathbb{R}^d$. Even if all model parame-
 96 ters were known exactly, the process y cannot be represented perfectly. This
 97 is due to unavoidable modelling assumptions, simplifications, or incomplete
 98 knowledge of the underlying physics. This disparity is known as *model dis-*
 99 *crepancy* [13] and is denoted by ϵ_{md} . The modelled physical process can
 100 therefore be described by $y = \eta(\mathbf{x}_c) + \epsilon_{md}$, where \mathbf{x}_c is an input configura-
 101 tion, such that $\eta(\mathbf{x}_c)$ summarizes all of the information the simulator carries
 102 about the system. Finally¹, the value of $\eta(\mathbf{x})$ is unknown until the model is
 103 evaluated at the input combination \mathbf{x} . When the model is computationally
 104 expensive, the analyst will only be able to run the model in a limited number
 105 of input configurations, which induces another source of uncertainty, called
 106 *code uncertainty* [14].

107 Given the sources of uncertainty introduced above, HM is designed to
 108 explore the input space \mathcal{X} and discard regions which are unlikely to produce
 109 the measured system response. This is achieved through: (i) the use of an
 110 *implausibility measure*, which quantifies the distance between the measure-
 111 ment z and the output of the model $\eta(\mathbf{x})$, normalized by the different sources
 112 of uncertainty; and (ii) a *Bayesian emulator* to alleviate the cost of running
 113 the complex model. The technical details behind Bayesian emulation are

¹If the simulator is stochastic in nature, *i.e.* evaluating η at a fixed input combination \mathbf{x} returns a different output value, $\eta(\mathbf{x})$ every time, another source of uncertainty called *ensemble variation* can be added. See discussion in Section 2.3

114 given in Section 2.2. Due to the use of an emulator, HM becomes an itera-
115 tive procedure in practice. At each iteration, also called *wave* [8], HM builds
116 increasingly accurate emulators and the implausibility measure provides a
117 rule to discard the subsets of the input domain that are unlikely to pro-
118 duce an acceptable match between model output and observed data. Once
119 the non-implausible domain in the current wave is identified, HM selects a
120 handful of points at which to evaluate $\eta(\cdot)$. This data is then used to refine
121 the approximation provided by the emulator in the non-implausible domain.
122 The process continues until a predefined stopping condition is satisfied.

123 In contrast to conventional calibration methods, which seek a single point
124 \mathbf{x}_c , HM identifies a set of input combinations that are likely to produce
125 a match between model prediction and measured data, within some level of
126 uncertainty. Furthermore, whereas standard Bayesian calibration will always
127 find a posterior distribution for acceptable inputs, HM can discover that the
128 model is an inadequate representation of the physical process by returning
129 an empty non-implausible domain. Thus, some authors regard HM as a
130 pre-calibration strategy [8].

131 Generating an initial design to run the simulator is the first step for a
132 typical HM workflow. Initially, the whole input domain \mathcal{X} is considered. To
133 explore the model output across the input domain with as few data points
134 as possible, a design with good space-filling qualities is generated. A Latin
135 hypercube sampling (LHS) plan [15] is often used. As suggested in [16], a
136 common choice is to have the number of sample points equal to $n = 10d$,
137 where d is the dimension of the input. In practice, the choice of n is often de-
138 termined by the computational budget. Once an LHS design is specified, the

139 simulator $\eta(\cdot)$ is evaluated at each input point \mathbf{x}_i , producing a corresponding
140 output $\eta(\mathbf{x}_i)$ for $i = 1 \dots n$. The resulting input-output pairs constitute an
141 experimental design denoted by $\mathcal{D} = \{\mathbf{x}_i, \eta(\mathbf{x}_i)\}_{i=1}^n$.

142 2.2. Emulation

143 An emulator is a statistical approximation to the output of an expensive
144 computer model. Gaussian process emulators [17] are widely used to infer
145 the output of expensive simulators based on a small number of training runs.
146 In this case, the experimental design \mathcal{D} defined in the previous subsection
147 provides such runs. Emulators provide a full probabilistic characterisation
148 of the output at untried input configurations. Their widespread use is due
149 to the fact that they not only provide a fast surrogate to the output of
150 the simulator, but also produce an analytic expression for the uncertainty
151 arising due to the limited number of model evaluations (referred to in Section
152 2.1 as code uncertainty). The applications of Gaussian process emulators
153 span many fields of science and technology. Some recent examples include
154 modelling submarine sliding and tsunami formation, [18] and reducing the
155 cost of engineering reliability analysis [19].

156 It is important to note that the original HM approach presented in [20]
157 and [21] is based on the concepts of Bayes linear emulation [22], which uses
158 mathematical expectation, instead of probability, as a primitive. This further
159 aids the mitigation of computational cost. In this paper however, we assume
160 a Bayesian emulator is of the form:

$$\hat{\eta}(\mathbf{x}) = h(\mathbf{x})^\top \boldsymbol{\beta} + Z(\mathbf{x}) \quad (1)$$

161 where $h : \mathbb{R}^d \rightarrow \mathbb{R}^q$ is a vector of known functions, $\boldsymbol{\beta} \in \mathbb{R}^q$ is a vector
162 of coefficients and $Z(\mathbf{x})$ is a zero mean Gaussian process with covariance
163 function $\sigma^2 c(\mathbf{x}, \mathbf{x}'; \boldsymbol{\psi})$, also known as covariance kernel [23]. The regression
164 term $h(\mathbf{x})^\top \boldsymbol{\beta}$ models the global trend of the output, whilst the Gaussian
165 process models local variations. The covariance of the Gaussian process at
166 two distinct inputs, \mathbf{x} and \mathbf{x}' , is the product of a process variance parameter
167 σ^2 and a positive semi-definite correlation function $c(\cdot, \cdot; \boldsymbol{\psi})$, parameterised
168 by $\boldsymbol{\psi}$. In this work, the Matérn (5/2) correlation function [24] is employed.
169 This function was chosen because it is stationary and because it exhibits
170 a moderate degree of smoothness, which is suitable for many applications
171 [25]. The Matérn (5/2) correlation function has the following mathematical
172 expression:

$$c(\mathbf{x}, \mathbf{x}'; \boldsymbol{\psi}) = \left(1 + \frac{\sqrt{5}\delta(\mathbf{x}, \mathbf{x}')}{\boldsymbol{\psi}} + \frac{5\delta^2(\mathbf{x}, \mathbf{x}')}{3\boldsymbol{\psi}^2} \right) \exp \left(-\frac{\sqrt{5}\delta(\mathbf{x}, \mathbf{x}')}{\boldsymbol{\psi}} \right) \quad (2)$$

173 where $\delta(\mathbf{x}, \mathbf{x}')$ is the Euclidean distance between \mathbf{x} and \mathbf{x}' .

174 In order to estimate the values of each of the parameters $\boldsymbol{\beta}$, σ^2 and $\boldsymbol{\psi}$,
175 prior probability distributions can be imposed, and their posterior distribu-
176 tions can be computed by conditioning on the training runs \mathcal{D} . A weak prior
177 [26] can be used for $\boldsymbol{\beta}$ and σ^2 , namely

$$p(\boldsymbol{\beta}, \sigma^2) \propto \sigma^{-2} \quad (3)$$

178 Conditional on \mathcal{D} , the two parameters are distributed according to a normal-

179 inverse-gamma distribution [27], with expected values given by

$$\hat{\boldsymbol{\beta}} = (\mathbf{H}^\top \mathbf{C}^{-1} \mathbf{H})^{-1} \mathbf{H}^\top \mathbf{C}^{-1} \mathbf{f} \quad (4)$$

$$\hat{\sigma}^2 = \frac{\mathbf{f}^\top (\mathbf{C}^{-1} - \mathbf{C}^{-1} \mathbf{H} (\mathbf{H}^\top \mathbf{C}^{-1} \mathbf{H}^{-1}) \mathbf{H}^\top \mathbf{C}^{-1}) \mathbf{f}}{n - q - 2} \quad (5)$$

180 where $\mathbf{H} = [h(\mathbf{x}_1), \dots, h(\mathbf{x}_n)]^\top$, $\mathbf{C} \in \mathbb{R}^{n \times n}$ such that $C_{ij} = c(\mathbf{x}_i, \mathbf{x}_j; \boldsymbol{\psi})$, and
 181 $\mathbf{f} = [\eta(\mathbf{x}_1), \dots, \eta(\mathbf{x}_n)]^\top$. The posterior distribution of $\boldsymbol{\psi}$ can be computed
 182 using a full Bayesian approach [28, 29]. Due to the potentially high computa-
 183 tional cost of Bayesian computations, some authors instead prefer resorting
 184 to maximum likelihood estimation [30].

185 It can be shown [17] that, conditional on the parameter estimates in
 186 Eq. (4) and Eq. (5), the posterior predictive distribution for the simulator
 187 output is

$$\eta(\mathbf{x}) \sim m(\mathbf{x}) + \sigma_c(\mathbf{x}) t_{n-q} \quad (6)$$

188 where t_{n-q} is the Student's-t distribution with $n - q$ degrees of freedom. The
 189 emulator's posterior mean $m(\mathbf{x})$ and posterior variance $\sigma_c^2(\mathbf{x})$ are given by

$$m(\mathbf{x}) = h(\mathbf{x})^\top \hat{\boldsymbol{\beta}} + t(\mathbf{x})^\top \mathbf{C}^{-1} (\mathbf{f} - \mathbf{H} \hat{\boldsymbol{\beta}}) \quad (7)$$

$$\begin{aligned} \sigma_c^2(\mathbf{x}) = & \hat{\sigma}^2 [c(\mathbf{x}, \mathbf{x}) - t(\mathbf{x})^\top \mathbf{C}^{-1} t(\mathbf{x}) \\ & + (h(\mathbf{x})^\top - t(\mathbf{x})^\top \mathbf{C}^{-1} \mathbf{H}) (\mathbf{H}^\top \mathbf{C}^{-1} \mathbf{H})^{-1} \\ & \times (h(\mathbf{x})^\top - t(\mathbf{x})^\top \mathbf{C}^{-1} \mathbf{H})^\top] \end{aligned} \quad (8)$$

190 where $t(\mathbf{x}) = [c(\mathbf{x}, \mathbf{x}_1; \boldsymbol{\psi}), \dots, c(\mathbf{x}, \mathbf{x}_n; \boldsymbol{\psi})]^\top$.

191 *2.3. Implausibility threshold*

192 Let $I : \mathbb{R}^d \rightarrow \mathbb{R}$ be a function that measures the implausibility that an
 193 input configuration \mathbf{x} will produce a simulator output matching the experi-
 194 mental observation z . When the simulator is expensive, this implausibility
 195 can be defined as the distance between z and the emulator mean, $m(\mathbf{x})$.
 196 This distance can be normalised in order to express it in terms of the num-
 197 ber of standard deviations of the overall uncertainty [1]. This results in the
 198 following expression:

$$I(\mathbf{x}) = \frac{|z - m(\mathbf{x})|}{\sqrt{\sigma_{me}^2 + \sigma_{md}^2 + \sigma_c^2(\mathbf{x})}} \quad (9)$$

199 where σ_{me}^2 and σ_{md}^2 are respectively the variances of the measurement error
 200 term ϵ_{me} and the model discrepancy term ϵ_{md} , as defined in Section 2.1.
 201 The term $\sigma_c^2(\mathbf{x})$, corresponding to the (current level) emulator’s posterior
 202 predictive variance in Eq. (8), quantifies the code uncertainty. In this work,
 203 the simulator is assumed deterministic: for the same input configuration, the
 204 output is fixed. It is possible to add a term in the denominator of equation
 205 Eq. (9) to account for *ensemble variability* in case the simulator is stochastic
 206 [1].

207 Suppose that the implausibility measure $I(\cdot)$ is evaluated at a particular
 208 sample point \mathbf{x}^* . For $I(\cdot)$ to be meaningful, it should be true that the smaller
 209 the value of $I(\mathbf{x}^*)$, the more likely it is that \mathbf{x}^* yields an output that matches
 210 the experimental data within the specified level of uncertainty. A criterion
 211 for setting an implausibility threshold is provided by Pukelsheim’s rule [31],
 212 which states that if a random variable X has a unimodal distribution with

213 mean μ and standard deviation σ , such as the Student's-t in Eq. (6), then

$$P(|X - \mu| > 3\sigma) \leq 0.05 \quad (10)$$

214 Hence, a natural criterion for accepting \mathbf{x}^* as a non-implausible input
215 combination is $I(\mathbf{x}^*) \leq 3$. Sample points that fail this criterion are con-
216 sidered implausible. The new wave of HM begins by sampling from the
217 non-implausible domain identified using this rule.

218 *2.4. Sampling design for new waves*

219 The initial design \mathcal{D} to train the emulator can be generated through
220 LHS. After the first wave, sampling from the non-implausible domain can
221 become challenging very quickly. This can be due to, for example, rapid
222 reduction in its size. An additional challenge is that the non-implausible
223 domain may become disconnected or exhibit a complex topology, which can
224 further complicate the sampling procedure.

225 The most intuitive strategy to deal with this problem is to generate an
226 LHS plan on the whole input space \mathcal{X} , then discard all implausible points,
227 determined by $I(\mathbf{x})$. This simple acceptance-rejection strategy can quickly
228 become inefficient if the non-implausible domain reduces to a small frac-
229 tion of \mathcal{X} . Multiple solutions have been proposed in the literature. An
230 implausibility-driven evolutionary Monte Carlo algorithm (IDEMC) was pro-
231 posed in [32]. This generates uniform designs for the target space using an
232 implausibility ladder, which might be challenging to determine. Another ap-
233 proach, discussed in [1] is to generate normally distributed samples centered
234 on each point from the non-implausible domain of the current wave. The

235 covariance matrix of the non-implausible samples is scaled to give a rela-
236 tively flat distribution. The challenge in this approach is to determine an
237 optimal scaling factor, which determines the rate at which the input space
238 is explored. As an alternative to the above methods, this paper proposes
239 to sample the non-implausible domain using subset simulation, a rare-event
240 sampling method used in engineering reliability analysis.

241 3. Subset simulation

242 Subset simulation (SuS) is an advanced Monte Carlo method that ef-
243 ficiently estimates probabilities of failure of engineering systems [33]. Let
244 $g : \mathbb{R}^d \rightarrow \mathbb{R}$ be a *performance function* used to model a physical system.
245 That is, $g(\cdot)$ encodes all the available information about the system's be-
246 haviour and attributes, such as its geometry, material properties and loads.
247 When the system is large and complex, specifying deterministic inputs of the
248 performance function can be unrealistic. Thus, the inputs \mathbf{x} can be modelled
249 as distributed according to a joint probability density function (PDF) $\pi(\mathbf{x})$ ².
250 The output of $g(\cdot)$ then becomes a random variable $Y = g(\mathbf{x})$, and failure
251 is formulated as the exceedance of this random variable over a prescribed
252 threshold $b \in \mathbb{R}$. The main interest of reliability analysis is to determine the
253 *probability of failure* $P(Y > b)$, given by

²Even though precise characterisation for the inputs can be specified, one may want to investigate different scenarios by varying those inputs according to some probability distributions

$$P_F = P(Y > b) = \int \pi(\mathbf{x}) \mathbf{1}(\mathbf{x} \in \mathcal{F}) d\mathbf{x}, \quad (11)$$

where F denotes the *failure event* defined as

$$F = \{Y > b\} = \{\mathbf{x} \in \mathcal{F}\} = \{\mathbf{x} \in \mathbb{R}^d : g(\mathbf{x}) > b\} \quad (12)$$

254 and $\mathcal{F} \subseteq \mathbb{R}^d$ is the failure region of the input space. The indicator function
 255 $\mathbf{1}(\cdot)$ is equal to 1 if $\mathbf{x} \in \mathcal{F}$ and is zero otherwise.

256 The idea behind SuS is to model F as contained in a nested sequence
 257 of events $F = F_m \subset F_{m-1} \subset \dots \subset F_1 \subset F_0 = \{\mathbf{x} \in \mathcal{X}\}$ such that the
 258 probability of failure can be computed as

$$P_F = P\left(\bigcap_{i=1}^m F_i\right) = P(F_1) \times P(F_2|F_1) \times \dots \times P(F_m|F_{m-1}) \quad (13)$$

259 This means that sampling from F is done by sampling progressively from
 260 more frequent conditional events. Every intermediate failure event corre-
 261 sponds to an iteration *level* in the SuS algorithm, whereby level 0 corresponds
 262 to initial Monte Carlo sampling of the whole input space \mathcal{X} . There are two
 263 important parameters in the algorithm: the level probability, denoted by p_0
 264 and defined as $p_0 \equiv P(F_i|F_{i-1})$, and the number of samples in each level, N .
 265 Both are determined by the modeller, such that p_0N and $1/p_0$ are integers.
 266 The level probability p_0 directly influences the properties of the estimator
 267 for P_F . The recommended range to minimise its coefficient of variation is
 268 $p_0 \in [0.1, 0.3]$ [34]. The number of samples at each level, N , can be set
 269 to achieve a given coefficient of variation in the estimation of P_F . In our
 270 experience, however, its prescribed value is mainly driven by the available

271 computational budget. It is worth noting that, in industrial settings, the
 272 performance function $g(\cdot)$ is rarely analytical or inexpensive to compute. In
 273 practice, it usually consists of one or more expensive computer model. Differ-
 274 ent authors (see [19] and references therein) have proposed different strategies
 275 to tackle this cost, some of which include using the emulators discussed in
 276 Section 2.2.

277 The SuS algorithm proceeds as follows. At the unconditional level 0,
 278 SuS starts by generating N independent samples $\mathbf{x}_1, \dots, \mathbf{x}_N \sim \pi(\mathbf{x})$. The
 279 performance function $g(\cdot)$ is evaluated and the corresponding output values
 280 are sorted in descending order, resulting in the list $\{b_k^{(0)} : k = 1, \dots, N\}$. The
 281 value $b_k^{(0)}$ gives the estimated output value corresponding to the exceedance
 282 probability $p_k^{(0)} = P(Y > b_k^{(0)})$ where

$$p_k^{(0)} = \frac{k}{N}, \quad k = 1, \dots, N. \quad (14)$$

283 The first intermediate failure level, b_1 is defined as the midpoint between
 284 $b_{p_0N}^{(0)}$ and $b_{p_0N+1}^{(0)}$. This way, the conditional failure relation

$$p_0 = P(Y > b_1) = P(F_1|F_0), \quad (15)$$

285 is satisfied. Note that, by construction, the p_0N top-ranked samples have
 286 responses greater or equal to b_1 . Thus, they are guaranteed to belong to
 287 the first intermediate failure level \mathcal{F}_1 . The generation of new samples from
 288 \mathcal{F}_1 is done by exploiting this property. The p_0N top-ranked samples are
 289 used as seeds to generate independent Markov chains from the target density
 290 $\pi(\mathbf{x}|F_1) \propto \pi(\mathbf{x})\mathbf{1}(\mathbf{x} \in \mathcal{F}_1)$. This results in generating $N_c = p_0N$ Markov

291 chains, each with length

$$N_s = \frac{N}{N_c} = \frac{1}{p_0} \quad (16)$$

292 Since the seeds already belong to the intermediate failure domain \mathcal{F}_1 , there is
 293 no burn-in period, usually required in MCMC simulations to generate a single
 294 Markov chain. The MCMC scheme employed in the original SuS algorithm
 295 [12], was the *modified Metropolis algorithm*, which uses a component-wise
 296 Metropolis-Hastings sampling to generate the Markov chains. Throughout
 297 the years, different strategies have been proposed and developed. An account
 298 of those strategies can be consulted in [35].

299 Subset simulation follows the same principle iteratively: the i^{th} level (for
 300 $i \geq 1$) is defined as $F_i = \{Y > b_i\}$, where b_i is determined as the midpoint
 301 between $b_{N_c}^{(i-1)}$ and $b_{N_c+1}^{(i-1)}$. Thus, at each intermediate failure level, the equa-
 302 tion $p_0 = P(F_i|F_{i-1})$ is satisfied. At level i , N_c independent Markov chains
 303 are generated from the target density $\pi(\cdot|F_i)$, each with length N_s . The pro-
 304 cess is repeated until the target threshold level b is reached. As before, let
 305 m denote the final intermediate level. The threshold value satisfies $b_m \geq b$
 306 and thus the number of conditional samples with responses greater than b ,
 307 exceeds N_c . The estimate of the failure probability is derived from Eq. (13),
 308 which can be written as

$$\hat{P}_F = p_0^{m-1} \frac{1}{N} \sum_{k=1}^N \mathbf{1}(\mathbf{x}_k \in \mathcal{F}_m), \quad (17)$$

309 where $\frac{1}{N} \sum_{k=1}^N \mathbf{1}(\mathbf{x}_k \in \mathcal{F}_m)$ is the estimate of the conditional failure proba-
 310 bility at level m .

311 Subset simulation is capable of efficiently sampling from disconnected
 312 failure regions that are, potentially, orders of magnitude smaller than the
 313 original input space. In order to illustrate this, consider the performance
 314 function $g : [0, 1]^2 \rightarrow \mathbb{R}$ given by

$$g(\mathbf{x}) = \sum_{i=1}^9 w_i \phi(\mathbf{x} | \mu_i, \mathbf{C}_i) \quad (18)$$

315 where each $w_i \in (0, 1)$ is a weight, μ_i the mean and \mathbf{C}_i the covariance ma-
 316 trix of the i^{th} Gaussian PDF, $\phi(\mathbf{x} | \mu_i, \mathbf{C}_i)$. The numerical values of these
 317 parameters are given in Table A.1 and the level contours of $g(\mathbf{x})$ are shown
 318 in Figure 1. Let the failure threshold be $b = 2.75$. The failure domain would
 319 then be $\mathcal{F} = \{\mathbf{x} \in [0, 1]^2 : g(\mathbf{x}) > 2.75\}$, which results in the disjoint failure
 320 set shown in Figure 1(c). The successive subplots in Figure 1 depict how SuS
 321 steers the sampling towards \mathcal{F} .

322 [Figure 1 about here.]

323 The example above suggests that a natural analogy can be established be-
 324 tween the non-implausible domain introduced in Section 2 and a failure set.
 325 Firstly, both are defined by specifying a threshold (for non-implausibility and
 326 for failure, respectively). Secondly, both may be significantly smaller than
 327 the original input space. Thirdly, they may become disconnected. This mo-
 328 tivates treating the non-implausible domain within HM as if it were a failure
 329 set. The corresponding sampling can therefore be done using SuS within
 330 HM.

331 4. Proposed Approach

332 As discussed in the previous section, the main aim of SuS is to estimate
333 the probability of failure given by Eq. (11). In order to do so, the algorithm
334 produces samples within each intermediate failure domain and eventually
335 from the failure domain \mathcal{F} . The proposed approach for HM takes advantage
336 of this property, since the prime objective is to eventually sample from the
337 non-implausible domain. It should be noted that SuS has previously been
338 used in the context of calibration, for the estimation of parameter posterior
339 distributions [36, 37]. However, as discussed in Section 2.1, whilst Bayesian
340 calibration always delivers a posterior distribution, HM might determine that
341 the non-implausible domain is empty.

In order to use SuS within HM, sampling is done by treating the non-implausible domain as if it were the failure domain defined by

$$\mathcal{F} = \{\mathbf{x} : I(\mathbf{x}) < 3\} \quad (19)$$

342 where the implausibility measure $I(\mathbf{x})$, defined in Eq. (9), takes the role of
343 the performance function.

344 The proposed approach begins with sampling the input domain of the
345 computer model and evaluating it to get an initial data set, \mathcal{D}_1 . This data
346 set is split and then used to train and validate the initial GPE³ [27]. At
347 this point, the parameters of SuS are set as per Section 3. It is important
348 to note that the direction of the inequality in $\mathcal{F} = \{\mathbf{x} : I(\mathbf{x}) < 3\}$ is the
349 opposite to that of the inequality in the definition of the failure domain

³If the code is very computationally expensive, the emulator can be validated using cross-validation instead of a separate validation set (see Section 2.1 in [30]).

350 given in Eq. (12). This feature is accounted for by sorting the negative of the
 351 values of the implausibility function evaluated at the candidate samples. The
 352 algorithm progresses by sequentially discarding regions of the input domain,
 353 according to their implausibility $I(\mathbf{x})$, as explained in Section 2.3. When
 354 SuS converges, it returns the set \mathbf{X}_{SuS} , which belongs to the non-implausible
 355 domain. A subset of these samples, which maximises the predictive variance
 356 of the GPE at the current level is selected and denoted as \mathbf{X}_{add} . Other
 357 approaches to selecting \mathbf{X}_{add} exist, such as the maximin strategy outlined in
 358 [1]. The GPE for the ℓ^{th} wave of HM is trained using an augmented data set,
 359 $\mathcal{D}_\ell = \mathcal{D}_{\ell-1} \cup \{\mathbf{X}_{add}, \eta(\mathbf{X}_{add})\}$. It should be pointed out that, even if the model
 360 itself is a highly non-linear function, it becomes smoother in the plausible
 361 region as the latter shrinks after each level of HM, which in turn leads to an
 362 increase in the accuracy of the emulator [1]. At the same time, the training
 363 points become denser. The algorithm terminates once the code uncertainty,
 364 quantified through the emulator variance, becomes smaller than the other
 365 sources of error. The proposed approach is summarised in Algorithm 1.

366 The next two sections present applications of the proposed SuS-based
 367 HM. In both examples, the modified Metropolis algorithm is used to sample
 368 from the intermediate failure domains. This is not a constraint, since any
 369 of the MCMC sampling schemes reviewed in [35] could in principle be used.
 370 Previous work on the comparison of some of these schemes within SuS-based
 371 HM can be found in [38].

Algorithm 1 History matching with subset simulation

- 1: Provide an experimental observation z and relevant standard deviations:
 σ_{me} for observational uncertainty and σ_{md} for model discrepancy.
 - 2: Set parameter values for SuS: p_0, N .
 - 3: Define $\mathcal{F} = \{\mathbf{x} : I_\ell(\mathbf{x}) < 3\}$ where $I_\ell(\mathbf{x})$ is the implausibility measure for the emulator at ℓ^{th} wave of HM as defined in Eq. (9).
 - 4: Generate a space-filling plan, $\mathbf{X} \in \mathbb{R}^{n \times d}$ and form $\mathcal{D}_1 = \{\mathbf{x}_i, \eta(\mathbf{x}_i)\}_{i=1}^n$
 - 5: Train a GPE $\hat{\eta}_1(\mathbf{x}) \sim m_1(\mathbf{x}) + \sigma_{c_1}(\mathbf{x})t_{n-q}$ on \mathcal{D}_1 and validate it.
 - 6: $\ell \leftarrow 1$.
 - 7: **while** $\sigma_{me} < \sigma_{c_\ell}$ **and** $\sigma_{md} < \sigma_{c_\ell}$ **do**
 - 8: Sample from the non-implausible domain using SuS:
 - 9: **Subset simulation**
 - 10: Obtain an MC sample $\mathbf{X}_{SuS} \in \mathbb{R}^{N \times d} \sim \pi(\mathbf{x})$.
 - 11: $N_F \leftarrow 0$.
 - 12: $j \leftarrow 0$.
 - 13: **while** $N_F < p_0 N$ **do**
 - 14: $j = j + 1$.
 - 15: Evaluate $\hat{\eta}_\ell(\mathbf{X}_{SuS})$.
 - 16: Compute $I_{SuS} \equiv -I_\ell(\mathbf{X}_{SuS})$ and sort in descending order.
 - 17: Renumber \mathbf{X}_{SuS} to match the order of I_{SuS} .
 - 18: Select $\{\mathbf{x}_{SuS}^{(i)}\}_{i=1}^{p_0 N}$ as seeds for MCMC.
 - 19: Compute intermediate threshold $b_j = \frac{1}{2} \left[I_{SuS}^{(p_0 N)} + I_{SuS}^{(p_0 N + 1)} \right]$
 - 20: Define intermediate failure domain $F_j = \{I_{SuS} > b_j\}$
 - 21: Obtain a sample, \mathbf{X}_{SuS} , from $\pi(\mathbf{x}|F_j)$ using an MCMC scheme.
 - 22: $N_F = \sum_{i=1}^N \mathbf{1}(I_{SuS} > -3)$
 - 23: **end while**
 - 24: $\ell = \ell + 1$
 - 25: Let \mathbf{X}_{add} be a subset of points from \mathbf{X}_{SuS} .
 - 26: Construct $\mathcal{D}_\ell \leftarrow \mathcal{D}_{\ell-1} \cup \{\mathbf{X}_{add}, \eta(\mathbf{X}_{add})\}$
 - 27: Train a GPE $\hat{\eta}_\ell(\mathbf{x}) \sim m_\ell(\mathbf{x}) + \sigma_{c_\ell}(\mathbf{x})t_{n-q}$ on \mathcal{D}_ℓ .
 - 28: **end while**
-

372 **5. Calibration: wing weight reduction**

373 This section demonstrates the proposed approach by using HM to cal-
 374 ibrate a model of the weight of a light aircraft wing. Weight is a critical
 375 factor in aircraft design and ensuring the model at hand can reliably match
 376 experimental weights is of vital importance.

377 The analytical model considered here is derived from historical data and is
 378 given by

$$W = 0.036S_w^{0.758}W_{fw}^{0.0035} \left(\frac{A_w}{\cos^2 \Lambda} \right)^{0.6} q^{0.006} \lambda^{0.04} \left(\frac{100t_c}{\cos \Lambda} \right)^{-0.3} (N_z W_{dg})^{0.49} + S_w W_p \quad (20)$$

379 [Table 1 about here.]

380 Eq. (20) was introduced in its original form in [39]. The last term on the
 381 right hand side representing the weight of the paint on the wing was added
 382 in [30]. A brief description of the inputs of the model, together with their
 383 ranges is provided in Table 1.

384 A simulated observation for the wing weight was set at $z = 130\text{lb}$. A mea-
 385 surement error of $\pm 5\text{lb}$ was imposed, corresponding to a standard deviation
 386 of $\sigma_{me} = 1.7\text{lb}$. Since, in this case, z is a synthetic surrogate for a physical
 387 observation, there is no direct meaning to the term model discrepancy and
 388 it is identically 0. Despite this, if the observation were coming from a real
 389 physical measurement the discrepancy term would have had some nonzero
 390 value. For this example, the model discrepancy was set to $\sigma_{md} = 1$, a value
 391 sufficient to make sure σ_{md} is included in the procedure, yet small enough so
 392 as to not overpower the uncertainty coming from the simulated measurement.

393 The treatment of model discrepancy is an important problem in uncertainty
394 quantification and an area of research in itself, see for example [13]. Finally,
395 the simulator described by Eq. (20) is deterministic and has no ensemble
396 error.

397 Following the ideas outlined in Section 2.2 a Gaussian process emulator
398 was trained with 100 points from an LHS design. The global trend term in
399 Eq. (1) was chosen as $h(\mathbf{x}) = 1$ in this case, so that the Gaussian process
400 component of the emulator was responsible for taking into account any devia-
401 tions from the mean. This choice is subjective and was motivated by the lack
402 of knowledge of the general shape of the function. Specifying more complex
403 forms for $h(\mathbf{x})$ is possible and can be informed by exploratory analysis. The
404 samples in the training set were normalized in $[0, 1]$ due to the large variation
405 of the input scales. This preprocessing step facilitates the search for optimal
406 correlation lengths, $\boldsymbol{\psi}$, and makes the results more easily interpretable. A
407 genetic algorithm was used to search the likelihood of the emulator for $\boldsymbol{\psi}$,
408 while $\boldsymbol{\beta}$ and σ^2 were computed from the expected values of their respective
409 distributions, given in Eq. (4) and Eq. (5).

410 At each wave, SuS was used to sample the non-improbable domain with
411 6000 points per subset level, and each level was given a target probability,
412 $p_0 = 0.1$. Out of the final sample, 10 points from the non-improbable domain
413 were added to the design at locations where the predictive variance from
414 the emulator, given by Eq. (8), was the largest. The number of samples is
415 such that there is at least one point representing each input. Additionally,
416 sites with maximum predictive variance were chosen to rapidly reduce the
417 uncertainty about the non-improbable domain. The GPE was retrained after

418 each wave.

419 After 9 HM waves and 80 additional evaluations of the model in Eq. (20),
420 the standard deviation of the code uncertainty, σ_c had decreased below that
421 of the measurement error, σ_{me} and the analysis was terminated. Each wave
422 required between 5 and 6 SuS levels, implying that the probability of the
423 non-implausible domain is on the order of 10^{-6} to 10^{-5} . Figure 2 depicts the
424 *optical depth* projections of the input space, introduced by [8]. These projec-
425 tions show the logarithm, base 10, of the empirical probability of finding a
426 non-implausible sample in a given region of the input space, when projected
427 onto a two-dimensional subspace. In this manner, optical depth projections
428 provide a way to visualize the non-implausible domain conditioned on the
429 pair of inputs in each subfigure. To generate these plots, the input subspace
430 of each pair of inputs was discretised in a 20×20 grid of point values. The
431 remaining 8 dimensions, which vary between subfigures, were represented
432 by a 50,000 point LHS sample. In this manner, to produce a single opti-
433 cal depth plot, the emulator for the appropriate wave of HM was evaluated
434 $20 \times 20 \times 50,000 = 20,000,000$ times.

435 The panels in the lower and upper triangles of Figure 2 show the projec-
436 tion plots from the first and last wave of HM, respectively. Several obser-
437 vations can be made from these plots. Firstly, many of the two-dimensional
438 projections of the input space exhibit subtle, but quantifiable reduction in
439 area from the first to the last wave of the analysis. This behaviour can be
440 attributed to the function being relatively smooth and the fact that the GPE
441 mean was capable of representing it with reasonable accuracy early on in the
442 procedure. This is to say that even though the mean of the emulator was

443 able to match the non-implausible domain reasonably well, its distance from
444 the training sample caused the predictive variance of the GPE to be larger
445 than the other sources of uncertainty, preventing the analysis from terminat-
446 ing. For some projections, such as $\Lambda - q$ and $\lambda - q$ the whole space seems
447 to have been discarded as implausible. This outcome is due to the sample-
448 based nature of the optical depth plots and the difficulty of producing data in
449 the non-implausible domain, by uniformly sampling the hidden dimensions.
450 Secondly, the scale of the log-probabilities is indicative of the overall size
451 of the non-implausible domain, with between 1 and 230 samples per 50,000
452 producing acceptable matches. Simple-looking problems such as these, show
453 the inadequacy of rejection-based uniform sampling and emphasise the im-
454 portance of effective methods to identify the non-implausible domain in each
455 wave of HM. Finally, the plots reveal how *active* certain inputs are, which
456 could lead to better understanding of the underlying model. For example,
457 the quarter-chord sweep angle Λ is not identified as important for satisfying
458 the measurement z , as seen from the fact that non-implausible samples are
459 uniformly distributed in its range. Similar conclusions can be made for W_{fw} ,
460 q and W_p . On the other hand, the wing area, S_w , the aspect ratio, A_w and
461 the load factor, N_z in particular are all influential in producing an accept-
462 able match to a relatively light wing, confirming the engineering intuition
463 that smaller, less-loaded wings can be made lighter.

464 [Figure 2 about here.]

465 The effect of HM on the output of the model is shown in Figure 3. The
466 samples identified to belong to the final non-implausible domain results in

467 the output values shown in Figure 3(a). In this figure, there is a tendency
468 for the outputs to cluster to the upper boundary of the prescribed region.
469 This behaviour serves as an evidence to the restrictive target weight used
470 in the analysis. For comparison, one of the most recognizable general avi-
471 ation aircraft, Cessna 172, has a wing weight of approximately 236 lb [40].
472 Figure 3(b) depicts a kernel estimation of the final distribution of the wing
473 weights, which is considerably narrower than the one used to train the initial
474 GPE.

475 [Figure 3 about here.]

476 The correlation between samples from SuS can be calculated using the
477 procedure outlined in Section 6 of [12] to determine the quality of the in-
478 formation they provide. In the above example, the coefficient of variation,
479 accounting for sample correlation varies in $\delta = [0.039, 0.068]$.

480 To illustrate how SuS is capable of sampling more efficiently from the
481 non-implausible domain, HM for the wing weight model was repeated using
482 MC sampling instead of SuS. All other aspects of the analysis were kept the
483 same, except for the number of MC samples. Since MC extracts all of its
484 information in one step, as opposed to SuS, which uses levels, the number
485 of samples required by MC to explore \mathcal{X} is much larger. At each wave,
486 $n_{MC} = 324,000$ samples were generated in \mathcal{X} , out of which $m = 10$ samples
487 were to be added to the training set for the next wave emulator. In this
488 comparison, n_{MC} is calculated as the total number of samples used in HM
489 with SuS (9 waves, with 6 SuS levels each and 6000 samples per level). Due to
490 the small volume of the non-implausible domain at each wave, MC was unable
491 to populate it densely enough and as a result $m < 10$ samples were added in

492 each wave. This outcome reveals one of the important advantages of using
 493 SuS for sampling the non-implausible domain: unless the set of acceptable
 494 matches is empty, SuS is able to populate it according to requirement. The
 495 MC-based HM terminated after $\ell = 4$ waves, due to the inability of MC to
 496 find samples in the non-implausible domain. A total of 9 samples were added
 497 across the 4 waves, which gives a GPE equivalent to the one in Wave 2 in
 498 SuS-based HM. It must be pointed out that the efficiency of SuS compared
 499 to MC comes at the cost of producing samples that cannot be guaranteed
 500 to be uniformly distributed over the non-implausible domain. However, the
 501 coefficient of variation accounting for sample correlation δ can be computed,
 502 as it was done above. This allows the analyst to monitor efficiency. An
 503 interesting question arises when this coefficient of variation is unacceptably
 504 large, due to high sample correlation. A potential solution could involve
 505 designing a thinning strategy for the modified Metropolis algorithm, but this
 506 is the subject of future work.

507 Figure 4 shows the *minimum implausibility plots* for three pairs of inputs
 508 from MC-based HM in the top row, compared to the ones from SuS-based
 509 HM in the bottom row. These plots depict the minimum implausibility
 510 in the high-dimensional domain, if a given pair of inputs were fixed to a
 511 particular value [8]. These plots reveal several things. Firstly, for all pairs
 512 of inputs, the non-implausible domain differs in topology. A particularly
 513 noticeable difference exists in the space spanned by thickness-to-chord ratio
 514 t_c and wing fuel weight W_{fw} . History matching with MC sampling identifies
 515 the non-implausible domain to be much more diffuse than the one in the
 516 SuS-based analysis. That is, it is larger and at the same time has higher

517 implausibility overall. Secondly, the GPE trained on data with a dense set
518 of non-implausible samples from SuS achieves better accuracy compared to
519 the true function. In the case of the wing weight simulator, the code can be
520 run affordably without the need of an emulator, to explore the implausibility
521 landscape without code uncertainty. The implausibility threshold, $I(x) = 3$,
522 is shown on each of the contour plots as a black dashed line. In all three
523 cases, the agreement is better for the lower line of plots. Finally, it must
524 be noted that the efficiency of the HM process increases when using SuS
525 as a sampler, since the quality of the GPE in the non-implausible domain
526 increases more rapidly when using informative samples. This decreases both
527 the number of potentially costly code evaluations and the number of waves,
528 and thereby emulators, to be generated.

529 [Figure 4 about here.]

530 Note that, even if the model in this case study is not computationally
531 expensive, it demonstrates the challenges in calibrating models with even
532 moderately-sized input domains.

533 **6. Robust design: aircraft wing-engine matching**

534 The second application of SuS within HM presented here is *robust de-*
535 *sign*. In engineering, the term robust design refers to the process of seeking
536 not only an optimal mean value of a system performance metric, but also
537 to ensure that this optimum is insensitive to variations which could lead to
538 undesired system behaviour [41, 42]. The essence of the robust design prob-
539 lem is prescribing a target value for quantities of interest that determine the

540 performance of a system. The designer’s task is to then find one or more
541 design input configurations that deliver this target within certain tolerance.

542 Suppose that a target value for a quantity of interest is treated as if it
543 were an experimental measurement. Also, suppose that the corresponding
544 tolerance can be treated as the underlying uncertainties. This treatment
545 provides an analogy between matching a model output with experimental
546 data (given the sources of uncertainty) and matching a design target within
547 a prescribed tolerance. Therefore, the proposed SuS sampling for HM can
548 also be used to solve the robust design problem by identifying the set of
549 input values that yield an output consistent with a design target within
550 certain tolerance. This results in a reduced input space that can be further
551 explored by an analyst in the search for an optimal design. Since HM can
552 deliver an empty non-improbable domain, the designer might conclude that
553 there is no input configuration that complies with the system requirements,
554 given the current model. This information can be very valuable in terms of
555 improving the model or rethinking the feasibility of the design targets.

556 This section develops the idea with an application to aircraft design.
557 Subset simulation has previously been used in different optimisation-related
558 problems [43, 44]. However, to the authors’ knowledge, it has not been used
559 in robust design. The application proposed in this section demonstrates how
560 SuS-based HM can be used in contexts beyond model calibration.

561 *6.1. Problem description*

562 Modern aircraft are expected to operate within very stringent perfor-
563 mance and regulatory limits to reduce their environmental impact, whilst
564 keeping their profitability as a mode of transportation. Increasingly demand-

565 ing regulations are coming into effect worldwide, which impose bounds on the
566 amount of nitrous oxide (NO_x), among other greenhouse gases produced by
567 aircraft [45]. Such requirements necessitate a highly structured approach to
568 early stage aircraft design, acknowledging the complex nature of interactions
569 and dependencies between different systems. For the purposes of this study,
570 and following the work in [46], the conceptual aircraft is defined as a com-
571 bination of different wings and engines, in an approach known as *set-based*
572 *design*. Each wing and engine are in turn defined by the parameters given in
573 Table 2.

574 Whilst the modelling process is multi-disciplinary and multi-organisational,
575 here it is presented in an abstract form as a chain of coupled analyses im-
576 plemented in a tool called AirCADia [47]. AirCADia is a framework for
577 interactive composition and exploration of conceptual aircraft design config-
578 urations. In this case study, six parameters were varied within AirCADia
579 to achieve the target emissions value. In order to collect all required data,
580 the model was run on a Lenovo ThinkCentre M900 Tower, with an Intel[®]
581 Core[™] i7-6700, 3.4 GHz CPU. On this machine, each evaluation took 0.5
582 seconds.

583 [Table 2 about here.]

584 6.2. History matching NO_x

585 The level of NO_x emissions was selected as the target output variable
586 that would drive the design. Initially, a GPE was trained on $n = 60$ Latin
587 hypercube points, using a global trend term, $h(\mathbf{x}) = [1, \mathbf{x}]^\top$. The emula-
588 tor was validated with another $m = 40$ LHS samples to verify its accuracy

589 in representing AirCADia’s output. The plot of simulator outputs against
590 GPE predictions is displayed in Figure 5(a). It shows the degree to which
591 predictions from the emulator correspond with simulator observations. If the
592 GPE were a perfect predictor, the scatter would have lain along the 45 de-
593 gree dashed line. The error bars indicate the 95% credible interval associated
594 with each point. Most of the predicted points contain the 45 degree line in
595 their credible intervals. As seen in Eq. (6), each prediction from the emulator
596 follows a Student’s-t distribution. Therefore, the residuals between simulator
597 output and prediction should occupy the interval $[-2, 2]$ with around 95%
598 probability. These normalised residuals, often termed *individual prediction*
599 *errors* [27], are plotted against predictions in Figure 5(b). The residuals are
600 uniformly spread around 0 with no discernible patterns, or significant num-
601 ber of outliers. Jointly, these visual diagnostics suggest that the emulator is
602 a reasonably accurate representation of the simulator. After validation, the
603 test points were added to the design of experiments and the trend coefficients
604 in Eq. (4), and process variance in Eq. (5) were re-estimated.

605 [Figure 5 about here.]

606 After consultation with the developers of AirCADia, at the Centre for
607 Aeronautics at Cranfield University, the target range for NO_x was chosen
608 as 240 ± 10 lb over a 3000 nautical mile trip, including landing and take-off
609 [47]. The reader is reminded that the end goal of the robust design task
610 is to attain a pre-specified level, with tolerance, of a quantity of interest.
611 This is in contrast to the aim of optimisation, where, typically, the analyst
612 seeks to attain an optimum level of the quantity of interest, possibly subject
613 to constraints. The experts’ reasoning behind choosing the specific NO_x

614 target is not provided herein, since it is not in line with the main aim of
615 the paper. As outlined before, in the robust design setting the target range
616 can be treated as measurement plus corresponding uncertainties. Therefore,
617 all uncertainties for HM are accumulated into the measurement error term.
618 In order to ensure that the target range is respected, HM was carried out
619 with an error term which ensures that 95% of the responses will lie in the
620 correct region. Thus, the final values for the analysis were set as $z = 240$
621 and $\sigma_{me} = 3.33$.

622 In each wave of HM, SuS was run with $N = 6,000$ samples per level
623 and level probability, $p_0 = 0.1$. In the first wave, two levels of SuS were
624 required to populate the non-improbable domain, implying that its proba-
625 bility is on the order of 10^{-2} . The two levels sampled the emulator a total
626 of 12,000 times, obtaining over 3,500 samples in the non-improbable do-
627 main. For comparison, a direct Monte Carlo simulation would have required
628 approximately 350,000 samples on average to achieve a similar result. The
629 code uncertainty associated with some of the samples from SuS exceeded
630 the measurement error and therefore it was necessary to continue with the
631 analysis.

632 The analysis was terminated after three waves, when the emulator vari-
633 ance $\sigma_c^2(\mathbf{x})$ was reduced sufficiently in comparison with the imposed uncer-
634 tainty⁴. From the denominator in Eq. (9), it can be seen that in this example,
635 $\sigma_c^2(\mathbf{x})$ is the only source of uncertainty that is free to change. Once it be-

⁴Despite the seemingly quick running times of the simulator, the analyses would have taken approximately 5 hours for SuS and just over 2.5 days for DMC, if the simulator was sampled directly.

636 comes small in comparison with the other components, the implausibility
637 measure does not change significantly. Further substantial reduction in the
638 non-implausible domain becomes unlikely.

639 Figure 6 shows diagnostics from the final wave of HM. The sub-figures
640 in the upper triangle contain the optical depth plots and those in the lower
641 triangle show the minimum implausibility plots. Together, these two repre-
642 sentations visualize the extent of the non-implausible domain. In Figure 6,
643 it can be seen that the inputs relating to engine pressure (OPR, FPR) have
644 significant contribution to the value of NO_x , since their domain was signif-
645 icantly reduced to achieve the specified target range. In particular, it was
646 not likely to find matching outputs for high values of OPR and low values
647 of FPR, regardless of the values of the other inputs. An interesting interac-
648 tion is the one between sea-level static thrust (SLST) and wing aspect ratio
649 (A_W), which indicates that low powered engines must be matched to efficient,
650 slender wings to attain the required NO_x level.

651 [Figure 6 about here.]

652 The values of NO_x corresponding to the non-implausible samples are
653 shown in Figure 7. Note that the values of the emissions in Figure 7(a)
654 exceed the specified range. This is due to the code uncertainty introduced
655 using the emulator instead of the original code. This uncertainty can be
656 reduced further, but an increase in the computational cost of the analysis
657 will be incurred, owing to the additional code evaluations needed to refine
658 the surrogate model. Figure 7(b), provides a visual comparison between the
659 pre- and post-history matching distributions of the output.

660 [Figure 7 about here.]

661 7. Conclusions

662 A solution to an important problem in model calibration with history
663 matching was proposed. The solution involves the use of subset simulation
664 to generate samples from the non-implausible domain of an expensive com-
665 puter model. It was shown that, within history matching, the volume of
666 the non-implausible domain may shrink by several orders of magnitude as
667 compared to the original input space. Thus, the non-implausible domain was
668 treated as a failure set, analogous to that in engineering reliability analysis.
669 This allowed the use of subset simulation as an efficient sampler, which pro-
670 vided good coverage of the non-implausible domain with a moderate number
671 of samples. The method selected highly informative input configurations,
672 which were used to train a Bayesian emulator. This led to a reduction in
673 computational time and fast convergence of the analysis.

674 The advantages of the proposed approach were demonstrated in two ex-
675 amples. The first one dealt with the calibration of an analytical wing model
676 to match a restrictively low target weight. The second example showed how
677 the proposed approach can be used as a pre-processor for robust design in
678 an industrial context. Future research based on this work includes exploring
679 the link between history matching and robust design with several, possibly
680 conflicting, design objectives. Another problem that requires attention is
681 that of local variations in the behaviour of the simulator. This might require
682 fitting different emulators if the non-implausible domain is disconnected.

683 **Acknowledgements**

684 The authors are grateful for the kind support from Dr. Xin Chen, Dr.
685 Arturo Molina-Cristobal, Dr. Atif Riaz and Prof. Marin Guenov, from the
686 Centre for Aeronautics at Cranfield University. Their team develops and
687 maintains AirCADia, and provided the authors with access to the model.

688 F. A. DiazDelaO acknowledges the support of the Data-centric engineer-
689 ing programme at The Alan Turing Institute, where he was a visiting fellow
690 under the EPSRC grant EP/S001476/1.

691 **Author Contributions**

692 The method was conceived by F. A. DiazDelaO, Z. T. Gong and M. Beer.
693 Initial coding and experiments were carried out by Z. T. Gong. An earlier
694 version of the manuscript was written for a conference by Z. T. Gong, F. A.
695 DiazDelaO and M. Beer [48]. More detailed coding and experiments were
696 carried out by P. O. Hristov. An extensive revision of the manuscript was
697 carried out by P. O. Hristov and F. A. DiazDelaO.

698 **References**

- 699 [1] Andrianakis, I., Vernon, I.R., McCreesh, N., McKinley, T.J., Oakley,
700 J.E., Nsubuga, R., Goldstein, M., and White, R.G., Bayesian history
701 matching of complex infectious disease models using emulation: A tuto-
702 rial and a case study on HIV in Uganda., *PLoS Computational Biology*,
703 11(1):1–29, 2015.

- 704 [2] Steyerberg, E.W., Vickers, A.J., Cook, N.R., Gerds, T., Gonen, M.,
705 Obuchowski, N., Pencina, M.J., and Kattan, M.W., Assessing the Per-
706 formance of Prediction Models, *Epidemiology*, 21(1):128–138, 2010.
- 707 [3] Cheng, Q.B., Chen, X., Xu, C.Y., Reinhardt-Imjela, C., and Schulte, A.,
708 Improvement and comparison of likelihood functions for model calibra-
709 tion and parameter uncertainty analysis within a Markov chain Monte
710 Carlo scheme, *Journal of Hydrology*, 519:2202 – 2214, 2014.
- 711 [4] Gibson, G.J. and Renshaw, E., Estimating parameters in stochas-
712 tic compartmental models using markov chain methods, *Mathematical
713 Medicine and Biology: A Journal of the IMA*, 15(1):19–40, 1998.
- 714 [5] Alfi, M. and Hosseini, S.A., Integration of reservoir simulation, history
715 matching, and 4D seismic for CO₂-EOR and storage at Cranfield, Mis-
716 sissippi, USA., *Fuel*, 175:116 – 128, 2016.
- 717 [6] Anterion, F., History matching: A one day long competition: classical
718 approaches versus gradient method., In *First International Forum On
719 Reservoir Simulation*, 1998.
- 720 [7] O’Hagan, A., Bayesian analysis of computer code outputs: A tutorial.,
721 *Reliability Engineering & System Safety*, 91(10):1290 – 1300, 2006.
- 722 [8] Vernon, I., Goldstein, M., and Bower, R., Galaxy formation: A Bayesian
723 uncertainty analysis., *Bayesian Analysis*, 5(4):619–669, 2010.
- 724 [9] Craig, P.S., Goldstein, M., Rougier, J.C., and Seheult, A.H., Bayesian
725 forecasting for complex systems using computer simulators, *Journal of
726 the American Statistical Association*, 96:717–730, 2001.

- 727 [10] Williamson, D. and Blaker, A.T., Evolving Bayesian emulators for struc-
728 tured chaotic time series, with application to large climate models., *So-*
729 *ciety for Industrial and Applied Mathematics*, 2014.
- 730 [11] Vernon, I., Goldstein, M., and Bower, R., Galaxy formation: Bayesian
731 history matching for the observable universe, *Statistical Science*,
732 29(1):81–90, 2014.
- 733 [12] Au, S.K. and Beck, J., Estimation of small failure probabilities in high
734 dimensions by subset simulation., *Probabilistic Engineering Mechanics*,
735 16(4):263–277, 2001.
- 736 [13] Brynjarsdóttir, J. and O’Hagan, A., Learning about physical pa-
737 rameters: The importance of model discrepancy, *Inverse Problems*,
738 30(11):114007, 2014.
- 739 [14] Kennedy, M.C. and O’Hagan, A., Bayesian calibration of computer mod-
740 els, *Journal of the Royal Statistical Society: Series B (Statistical Method-*
741 *ology)*, 63(3):425–464, 2001.
- 742 [15] McKay, M.D., Beckman, R.J., and Conover, W.J., A comparison of
743 three methods for selecting values of input variables in the analysis of
744 output from a computer code., *Technometrics*, 21:239–245, 1979.
- 745 [16] Loepky, J.L., Sacks, J., and Welch, W.J., Choosing the sample size of
746 a computer experiment: A practical guide., *Technometrics*, 51:366–378,
747 2009.

- 748 [17] Oakley, J.E. and O'Hagan, A., Probabilistic sensitivity analysis of com-
749 plex models: A Bayesian approach., *Journal of the Royal Statistical*
750 *Society. Series B (Statistical Methodology)*, 66:751–770, 2004.
- 751 [18] Salmanidou, D.M., Guillas, S., Georgiopoulou, A., and Dias, F., Statis-
752 tical emulation of landslide-induced tsunamis at the Rockall Bank, N-E
753 atlantic, *Proceedings of the Royal Society A: Mathematical, Physical and*
754 *Engineering Sciences*, 473(2200):20170026, 2017.
- 755 [19] Hristov, P.O., DiazDelaO, F.A., Farooq, U., and Kubiak, K.J., Adap-
756 tive Gaussian process emulators for efficient reliability analysis, *Applied*
757 *Mathematical Modelling*, 71:138 – 151, 2019.
- 758 [20] Craig, P., Goldstein, M., Seheult, H., and Smith, J.A., Bayes lin-
759 ear strategies for matching for matching hydrocarbon reservoir history,
760 *Bayesian Statistics*, 5:69 – 95, 1996.
- 761 [21] Craig, P., Goldstein, M., Seheult, H., and Smith, J.A., Pressure match-
762 ing for hydrocarbon reservoirs: A case study in the use of Bayes linear
763 strategies for large computer experiments, In *Case Studies in Bayesian*
764 *Statistics*, pp. 37–93. Springer New York, 1997.
- 765 [22] Goldstein, M. and Wooff, D., *Bayes Linear Statistics, Theory and Meth-*
766 *ods*, Wiley Series in Probability and Statistics, Wiley, 2007.
- 767 [23] Rasmussen, C.E. and Williams, C.K.I., *Gaussian processes for machine*
768 *learning.*, MIT Press, 2006.
- 769 [24] Minasny, B. and McBratney, A.B., The Matérn function as a general
770 model for soil variograms., *Geoderma*, 128(3-4):192 – 207, 2005.

- 771 [25] Duvenaud, D., Automatic model construction with Gaussian processes,
772 PhD thesis, University of Cambridge, 2014.
- 773 [26] Oakley, J., Bayesian uncertainty analysis for complex computer codes.,
774 PhD thesis, University of Sheffield, 1999.
- 775 [27] Bastos, L.S. and O’Hagan, A., Diagnostics for Gaussian process emula-
776 tors., *Technometrics*, 51(4):425–438, 2009.
- 777 [28] Garbuno-Inigo, A., DiazDelaO, F.A., and Zuev, K.M., Transi-
778 tional annealed adaptive slice sampling for Gaussian process hyper-
779 parameter estimation, *International Journal for Uncertainty Quantifi-*
780 *cation*, 6(4):341–359, 2016.
- 781 [29] Garbuno-Inigo, A., DiazDelaO, F.A., and Zuev, K.M., Gaussian process
782 hyper-parameter estimation using parallel asymptotically independent
783 Markov sampling, *Computational Statistics & Data Analysis*, 103:367 –
784 383, 2016.
- 785 [30] Forrester, A.I.J., Sobester, A., and Keane, A.J., *Engineering Design Via*
786 *Surrogate Modelling: A Practical Guide*, J. Wiley, 2008.
- 787 [31] Pukelsheim, F., The three sigma rule., *The American Statistician*,
788 88:88–91, 1994.
- 789 [32] Williamson, D. and Vernon, I. Efficient uniform designs for multi-wave
790 computer experiments. arXiv:1309.3520, 2013.
- 791 [33] Au, S.K. and Patelli, E., Rare event simulation in finite-infinite dimen-

- 792 sional space., *Reliability Engineering and System Safety*, 148:67 – 77,
793 2016.
- 794 [34] Zuev, K.M., Beck, J.L., Au, S.K., and Katafygiotis, L., Bayesian post-
795 processor and other enhancements of subset simulation for estimating
796 failure probabilities in high dimensions, *Computers & Structures*, 92–
797 93:283–296, 2012.
- 798 [35] Papaioannou, I., Betz, W., Zwirgmaier, K., and Straub, D., MCMC
799 algorithms for subset simulation, *Probabilistic Engineering Mechanics*,
800 41:89 – 103, 2015.
- 801 [36] Chiachio, M., Beck, J.L., Chichio, J., and Rus, G., Approximate
802 Bayesian computation by subset simulation, *SIAM Journal of Scientific*
803 *Computing*, 36(3):A1339 – A1358, 2014.
- 804 [37] DiazDelaO, F., Garbuno-Inigo, A., Au, S., and Yoshida, I., Bayesian
805 updating and model class selection with subset simulation, *CMAME*,
806 317:1102–1121, 2017.
- 807 [38] Gong, Z.T., DiazDelaO, F.A., and Beer, M., Sampling schemes for his-
808 tory matching using subset simulation, In *Proceedings of the 2nd In-*
809 *ternational Conference on Uncertainty Quantification in Computational*
810 *Sciences and Engineering*, 2017.
- 811 [39] Raymer, D., *Aircraft Design: A Conceptual Approach*, AIAA Education
812 Series, American Institute of Aeronautics & Astronautics, 1999.
- 813 [40] Taylor, J.W. (Ed.), *Jane’s All the World Aircraft 1977-78*, Jane’s Infor-
814 mation Group, 68 edition, 1978.

- 815 [41] Beyer, H.G. and Sendhoff, B., Robust optimization – A comprehensive survey, *Computer methods in applied mechanics and engineering*,
816 196(33-34):3190–3218, 2007.
- 818 [42] Zang, C., Friswell, M., and Mottershead, J., A review of robust optimal
819 design and its application in dynamics, *Computers & structures*, 83(4-
820 5):315–326, 2005.
- 821 [43] Li, H.S. and Au, S.K., Design optimization using subset simulation algorithm., *Structural Safety*, 32(6):384 – 392, 2010.
- 823 [44] Dubourg, V., Sudret, B., and Bourinet, J.M., Reliability-based design
824 optimization using Kriging surrogates and subset simulation., *Structural
825 and Multidisciplinary Optimization*, 44(5):673–690, 2011.
- 826 [45] EASA, European Aviation Environmental Report, Tech. rep., European
827 Aviation Safety Agency, 2016.
- 828 [46] Riaz, A., Guenov, M.D., and Molina-Cristobal, A., Set-Based Approach
829 to Passenger Aircraft Family Design, *Journal of Aircraft*, 54(1):310–326,
830 2017.
- 831 [47] Guenov, M., Nunez, M., Molina-Cristóbal, A., Datta, V.C., and Riaz,
832 A., Aircadia – An Interactive Tool for the Composition and Exploration
833 of Aircraft Computational Studies At Early Design Stage, *29th Congress
834 of the International Council of the Aeronautical Sciences*, pp. 1–12, 2014.
- 835 [48] Gong, Z.T., DiazDelaO, F.A., and Beer, M., Bayesian model calibration
836 using subset simulation, In *Risk, Reliability and Safety: Innovating*

837 *Theory and Practice: Proceedings of ESREL 2016, Glasgow, Scotland.,*
838 *2016.*

839 **List of Figures**

840 1 Sampling from a small probability event via subset simulation.
841 (a) Samples from the unconditional failure domain \mathcal{F}_0 (i.e.
842 the entire input space); (b) samples in the first intermediate
843 failure domain $\mathcal{F}_1 \subseteq \mathcal{F}_0$; (c) samples in the failure domain
844 $\mathcal{F} \subseteq \mathcal{F}_1 \subseteq \mathcal{F}_0$ generated by MCMC. 44

845 2 Pairwise optical depth plots for the first (lower triangle) and fi-
846 nal (upper triangle) waves of history matching for wing weight.
847 The plots show the evolution of the size of the non-implausible
848 domain and reflect the decreasing log-probability of finding ac-
849 ceptable input combinations (color bar) in different regions of
850 the input space. (color online). 45

851 3 Wing weight realizations from the final wave of history match-
852 ing; (a) emulator predictions (blue dots) compared to the spec-
853 ified target range; (b) kernel density estimation of the initial
854 code output distribution (orange fill) and that from samples
855 in the final non-implausible domain (purple fill). Dashed lines
856 in both figures show the target range (color online). 46

857 4 Minimum implausibility plots for three pairs of inputs to the
858 wing weight model. Top row: wave 4 HM results with Monte
859 Carlo sampling. Bottom row: wave 2 HM results with SuS
860 sampling. Dotted line: decision boundary of the non-implausible
861 domain without code uncertainty (color online). 47

862 5 Predictive diagnostics for the NO_x GPE. (a) correlation be-
863 tween prediction and observations with 95% credible intervals
864 depicted as error bars; (b) individual prediction errors for the
865 validation set. 48

866 6 Pairwise minimum implausibility (lower triangle, left color
867 bar) and optical depth (upper triangle, right color bar) plots
868 from the last wave of NO_x history matching. The color bar
869 on the right depicts the empirical log-probability of finding a
870 non-implausible sample in a given region of the input domain,
871 whereas the one on the left indicates the expected implausi-
872 bility value of that sample. Inputs belonging to the “Engine”
873 subsystem are clearly affected more than those belonging to
874 the “Airframe” subsystem in Table 2 (color online). 49

875	7	History matching identifies input configurations, which result	
876		in output values lying in the specified target range (dashed	
877		lines); (a) emulator predictions (blue dots) and the observation	
878		error distribution; (b) kernel density estimation of code out-	
879		puts before (orange fill) and after (purple fill) history matching	
880		(color online).	50

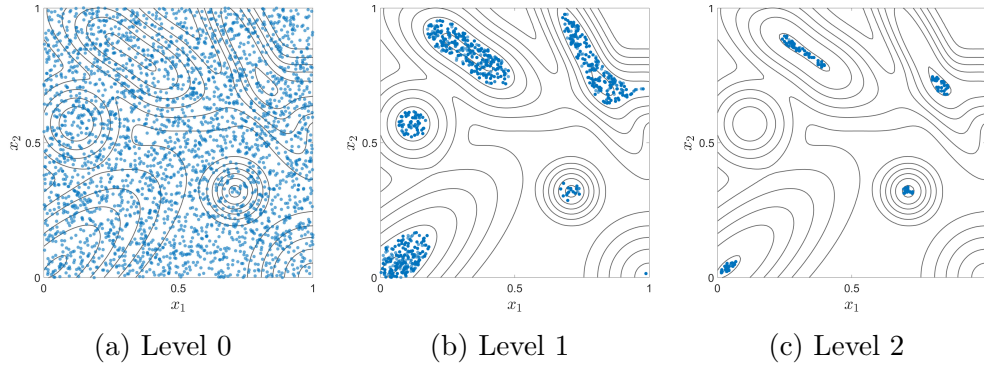


Figure 1: Sampling from a small probability event via subset simulation. (a) Samples from the unconditional failure domain \mathcal{F}_0 (i.e. the entire input space); (b) samples in the first intermediate failure domain $\mathcal{F}_1 \subseteq \mathcal{F}_0$; (c) samples in the failure domain $\mathcal{F} \subseteq \mathcal{F}_1 \subseteq \mathcal{F}_0$ generated by MCMC.

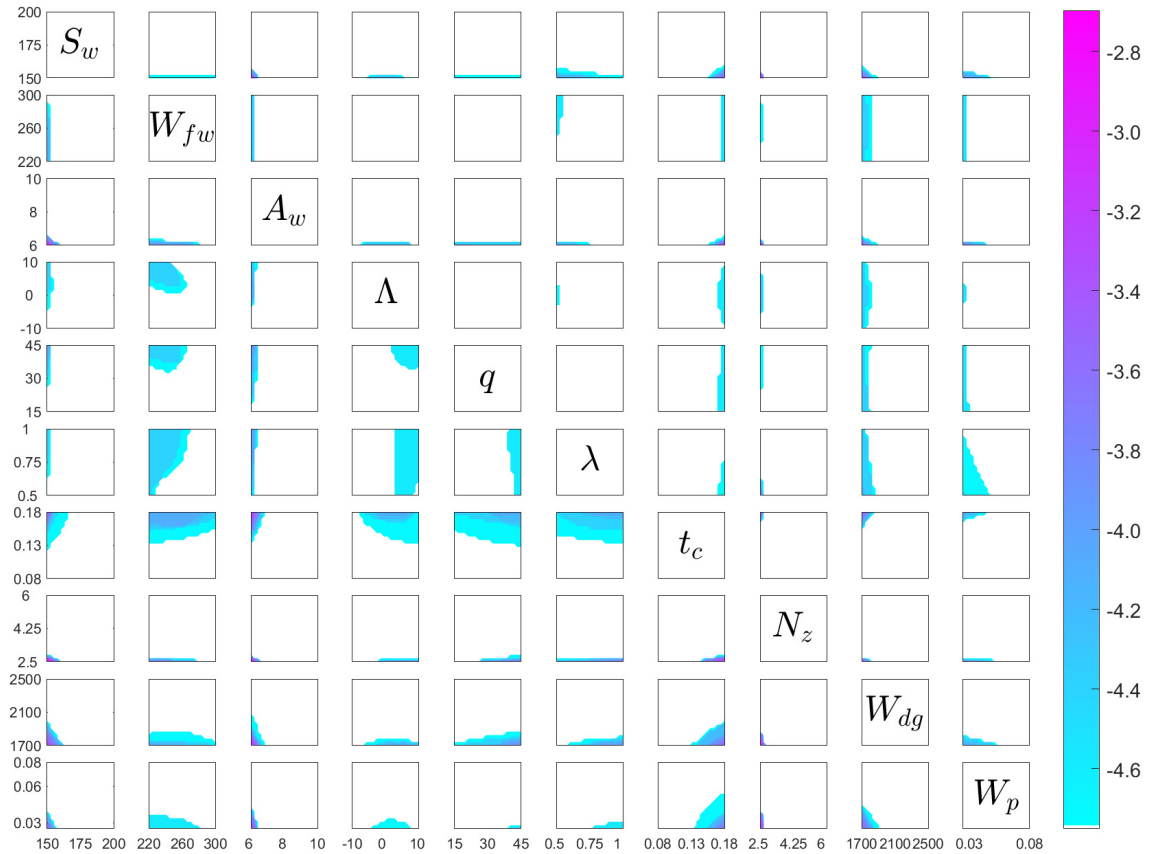


Figure 2: Pairwise optical depth plots for the first (lower triangle) and final (upper triangle) waves of history matching for wing weight. The plots show the evolution of the size of the non-implausible domain and reflect the decreasing log-probability of finding acceptable input combinations (color bar) in different regions of the input space. (color online).

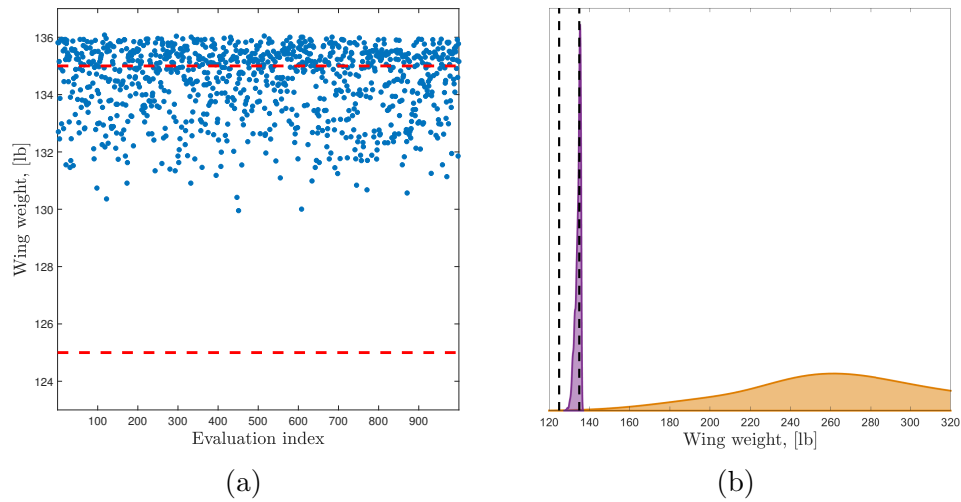


Figure 3: Wing weight realizations from the final wave of history matching; (a) emulator predictions (blue dots) compared to the specified target range; (b) kernel density estimation of the initial code output distribution (orange fill) and that from samples in the final non-implausible domain (purple fill). Dashed lines in both figures show the target range (color online).

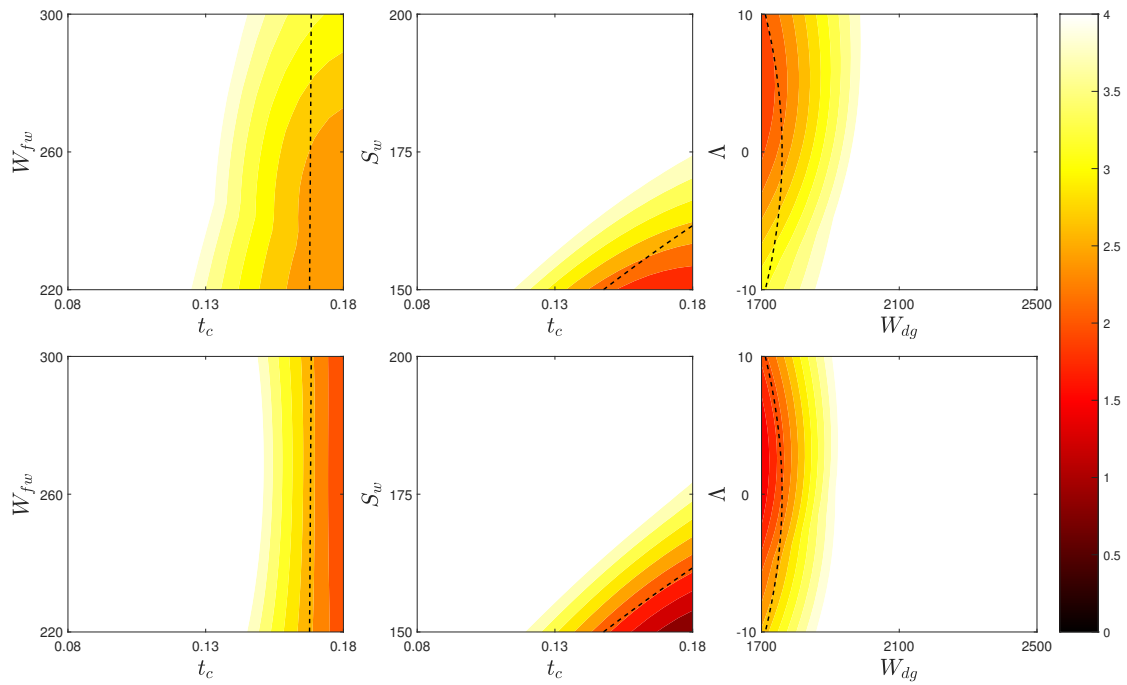


Figure 4: Minimum implausibility plots for three pairs of inputs to the wing weight model. Top row: wave 4 HM results with Monte Carlo sampling. Bottom row: wave 2 HM results with SuS sampling. Dotted line: decision boundary of the non-implausible domain without code uncertainty (color online).

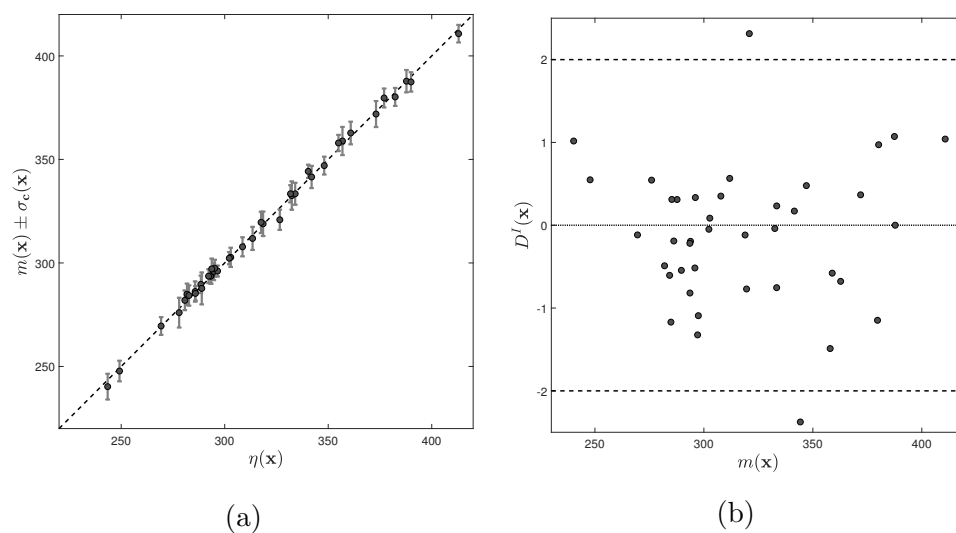


Figure 5: Predictive diagnostics for the NO_x GPE. (a) correlation between prediction and observations with 95% credible intervals depicted as error bars; (b) individual prediction errors for the validation set.

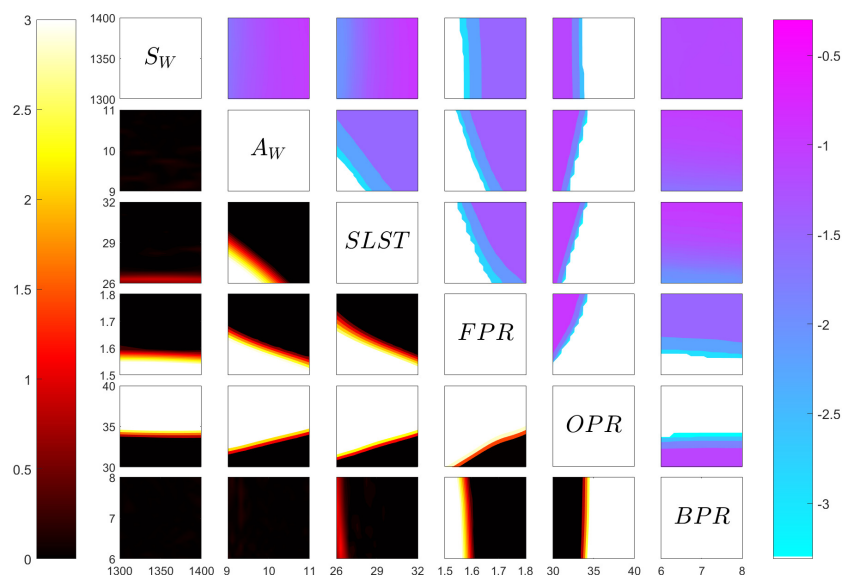


Figure 6: Pairwise minimum implausibility (lower triangle, left color bar) and optical depth (upper triangle, right color bar) plots from the last wave of NO_x history matching. The color bar on the right depicts the empirical log-probability of finding a non-implausible sample in a given region of the input domain, whereas the one on the left indicates the expected implausibility value of that sample. Inputs belonging to the “Engine” subsystem are clearly affected more than those belonging to the “Airframe” subsystem in Table 2 (color online).

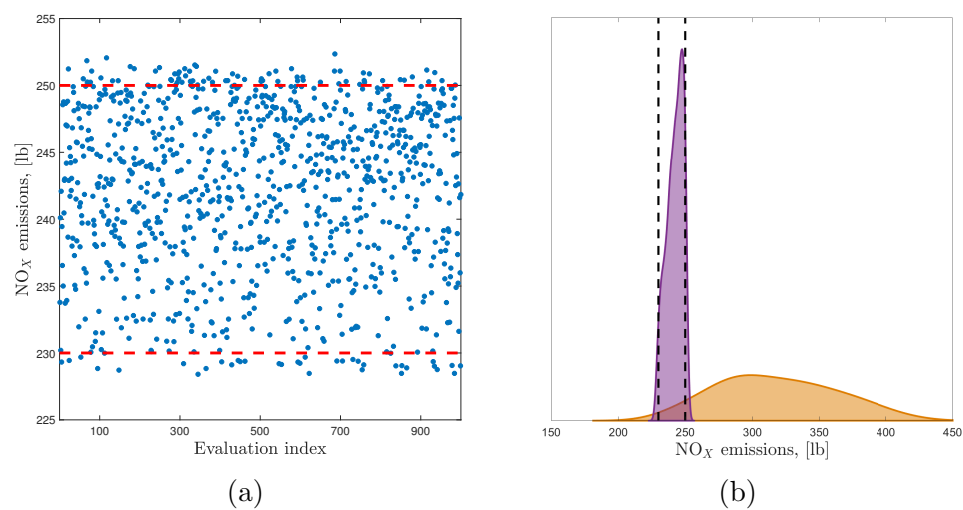


Figure 7: History matching identifies input configurations, which result in output values lying in the specified target range (dashed lines); (a) emulator predictions (blue dots) and the observation error distribution; (b) kernel density estimation of code outputs before (orange fill) and after (purple fill) history matching (color online).

881 **List of Tables**

882	1	Inputs for the light aircraft wing weight model.	52
883	2	Inputs and output for the climb-cruise case study, with respec-	
884		tive parent system and target ranges.	53
885	A.1	Parameter values for $g(\mathbf{x})$ in Eq. (18).	56

Table 1: Inputs for the light aircraft wing weight model.

Notation	Name	Range	Unit
S_w	Wing area	[150, 200]	ft ²
W_{fw}	Weight of fuel in the wing	[220, 300]	lb
A_w	Aspect ratio	[6, 10]	-
Λ	Quarter-chord sweep	[-10, 10]	deg
q	Dynamic pressure at cruise	[16, 45]	lb/ft ²
λ	Taper ratio	[0.5, 1]	-
t_c	Airfoil thickness to chord ratio	[0.08, 0.18]	-
N_z	Ultimate load factor	[2.5, 6]	-
W_{dg}	Design gross weight	[1700, 2500]	lb
W_p	Paint weight per unit area	[0.025, 0.08]	lb/ft ²

Table 2: Inputs and output for the climb-cruise case study, with respective parent system and target ranges.

Notation	Name	System	Range	Unit
S_W	Wing area	Airframe	[1300, 1400]	ft ²
A_W	Aspect ratio	Airframe	[9, 11]	-
SLST	Static thrust	Engine	[26, 32]	lbf $\times 10^3$
FPR	Fan pressure ratio	Engine	[1.5, 1.8]	-
OPR	Overall pressure ratio	Engine	[30, 40]	-
BPR	Bypass ratio	Engine	[6, 8]	-
NO _x	Nitrous oxide emissions	Output	240 \pm 10	lb

886 **Appendix A. SuS illustration function**

887 The function in Eq. (18), whose contour levels are shown in Figure 1 is a
888 mixture of nine bivariate Gaussian random variables with mean, covariance
889 and weight given in Table A.1.

890 [Table 3 about here.]

891 **List of Tables**

Table A.1: Parameter values for $g(\mathbf{x})$ in Eq. (18).

i	w_i	μ_i^\top	\mathbf{C}_i
1	0.327	[0.04 0.04]	$\begin{bmatrix} 0.030 & 0.020 \\ 0.020 & 0.025 \end{bmatrix}$
2	0.096	[0.98 0.70]	$\begin{bmatrix} 0.020 & 0 \\ 0 & 0.003 \end{bmatrix}$
3	0.143	[0.75 0.85]	$\begin{bmatrix} 0.010 & -0.015 \\ -0.015 & 0.030 \end{bmatrix}$
4	0.038	[0.71 0.32]	$\begin{bmatrix} 0.002 & 0 \\ 0 & 0.002 \end{bmatrix}$
5	0.161	[0.33 0.83]	$\begin{bmatrix} 0.020 & -0.010 \\ -0.010 & 0.010 \end{bmatrix}$
6	0.023	[0.43 0.73]	$\begin{bmatrix} 0.005 & 0 \\ 0 & 0.005 \end{bmatrix}$
7	0.026	[0.23 0.93]	$\begin{bmatrix} 0.005 & 0 \\ 0 & 0.005 \end{bmatrix}$
8	0.104	[1.00 0.00]	$\begin{bmatrix} 0.008 & 0 \\ 0 & 0.008 \end{bmatrix}$
9	0.081	[0.12 0.57]	$\begin{bmatrix} 0.005 & 0 \\ 0 & 0.005 \end{bmatrix}$