

1 Bayesian Uncertainty Quantification and Emulation of
2 Building Energy Models: a Tutorial

3 Dario Domingo^{a,*}, Mohammad Royapoor^b, Hailiang Du^a, Michael
4 Goldstein^a

^a*Department of Mathematical Sciences, Durham University, Stockton
Road, Durham, DH13LE, UK*
^b*Mohammad, please fill the affiliation, , , ,*

5 **Abstract**

Correct calibration of building energy models is crucial to ensure that these can be reliably used to make predictions and inform design decisions. This work provides guidance and methodological steps for a systematic treatment of uncertainty before and during calibration. We introduce emulation, a Bayesian technique to build a fast statistical surrogate of the model. The emulator's speed allows to explore the model's parameter space extensively in short timescales. The emulator's predictions, alongside several quantified uncertainties, are then used to sequentially rule out model parameter configurations which are highly implausible to replicate the observed data. This procedure is known as History Matching (HM) and has been successfully employed even in high dimensions and with slow simulators. A package for its implementation is referenced. We discuss and illustrate HM on an example of a building energy model, where both energy consumption and temperature data are matched. The results show that only 0.3% of the model input space leads to outputs that can match the observed energy consumption. The percentage reduces to 0.01% if temperature constraints are also accounted for. The example shows the efficiency of HM to locate the small region where model outputs and observations can match, and therefore its ability to deliver models that lead to more robust decision support.

6 *Keywords:* Uncertainty, Model Discrepancy, Building Energy Models,
7 History matching.

*Corresponding author

Email address: dario.domingo@durham.ac.uk (Dario Domingo)

8 **1. Introduction**

9 Up to 40% of primary energy use globally is attributed to buildings. Given
10 the global decarbonisation efforts, accurate characterisation of building en-
11 ergy has become a key tool to enable optimal design, retrofit and investment
12 decisions.

13 In order to produce high-quality modelling results that acknowledge mod-
14 elling limitations, full treatment of uncertainty becomes essential. This has
15 driven an accelerated trend in model uncertainty research and performance
16 gap analysis. A recent article has reported that a large proportion of the
17 building modelling community lacks essential knowledge on what the most
18 fundamental parameter inputs for buildings are and how these impact model
19 predictions [1]. This knowledge gap has major consequences, as retrofit
20 strategies are mostly derived in consultation with building energy models.
21 As such, techno-economic benefits of proposed solutions can be misleading.
22 In addition to this, future buildings are expected to become more respon-
23 sive to other civic activities (i.e. power generation and storage, transport,
24 etc.)[2]. Therefore imminent horizons require aggregated energy data at dis-
25 trict and potentially city level, where accuracy and treatment of uncertainty
26 are invaluable tools for decision makers at all levels.

27 A preliminary step to support model-based decisions is model calibration.
28 Calibration is the process through which values of the model parameters are
29 chosen, so that model outputs replicate observations within prescribed tol-
30 erances. Classically, this is performed by running the model multiple times,
31 for different values of the parameters. Hence, for each simulation, one or
32 more measures of the discrepancy between the corresponding output and the
33 observed data are evaluated. In the context of building energy consumption,
34 ASHRAE guidelines [3] are often followed [4, 5] to identify error thresholds
35 below which the discrepancy can be considered sufficiently small.

36 The main challenge of this approach is that the number of simulations
37 needed to cover uniformly the model’s parameter space grows exponentially
38 with increasing parameter dimensions. To meet the computational challenges
39 involved, a new area of Bayesian statistics has arisen, namely emulation [6–
40 9], part of the wider field of uncertainty quantification (UQ) of computer
41 models.

42 An emulator is a statistical surrogate of the simulator, providing pre-
43 dictions of the simulator’s outputs at inputs for which the original model
44 has never been evaluated. The emulator can be built and validated on a

45 personal device. It is instantaneous to run and therefore allows to predict
46 the model results for a very large number of inputs. Moreover, it quantifies
47 the uncertainty associated with such predictions. This explains the success of
48 emulation in a wide range of fields where complex computer models are used,
49 including but not limited to climate [10] and energy price projections [11],
50 epidemiology [12], hydrocarbon reservoirs [13] and galaxy formation [14].

51 Emulation is developed under Bayesian principles, arguably the most natural
52 framework to deal with uncertainties. In the context of calibration, key
53 uncertainties that are often overlooked are: i) the ability of the model to
54 simulate the process of interest (for example, building energy consumption);
55 ii) and the intrinsic accuracy of the observations used to calibrate the model.
56 To deal coherently with these and other sources of uncertainty, we propose
57 the use of so-called implausibility measures to compare model (or emulated)
58 outputs and observations.

59 Implausibility measures have been employed in the contexts of reconstruction
60 of ice sheet shapes [15], gene evaluation [16], galaxy formation [14] and
61 more, within a process known as History Matching (HM). The process sequentially
62 removes regions of the input space deemed implausible to produce
63 outputs which match the observations. However, in contrast to classical approaches
64 in building energy model calibration, it does so in light of quantified
65 uncertainties affecting various components of the system.

66 This work is aimed at researchers in the energy community, who are keen
67 to make their model-based inference more robust, by operating in a framework
68 that recognises and accounts for unavoidable sources of uncertainty.
69 The statistical framework is illustrated on an example study of a single-house
70 energy model, where model parameters are sought that simultaneously match
71 energy and environmental features. A link to an R package is provided, which
72 allows the researcher to implement the proposed methodology within their
73 own research.

74 This work is organised as follows. Section 2 discusses key sources of uncertainty
75 in computer model-based inference, particularly in the context of building energy
76 models. Limitations of current approaches in this field are also outlined.
77 Section 3 illustrates the proposed methodology: from emulation to
78 implausibility measures and HM. Section 4 gives details of the building
79 energy model and field data used in our case study. Results on the example
80 study are illustrated in Section 5, where we also discuss approaches to scenarios
81 that may be encountered in other case studies. Section 6 concludes
82 the article with discussions and work extensions.

83 **2. Uncertainty Sources in Building Energy Models**

84 *2.1. Calibration under Uncertainty*

85 Calibration is a crucial topic not just in energy, but in a variety of sci-
86 entific disciplines. It refers to the process by which the parameters of a
87 computer model are identified so that model outputs match observed data.
88 Mathematically, the model (or simulator) can be seen as a function f , with
89 input \mathbf{x} containing the list of model parameters of interest, and returning
90 $f(\mathbf{x})$ as output. If we call z the observed data we want to match, the cali-
91 bration task can be phrased as finding the input(s) \mathbf{x} so that $f(\mathbf{x})$ is “close
92 enough” to z . In practical applications, we often have a sequence of obser-
93 vations to match (and correspondingly a sequence of model outputs), rather
94 than a single observation.

95 Proximity between simulations and observations should be assessed upon
96 consideration of all sources of uncertainty that affect the system. While
97 an exhaustive list of such sources is problem-dependent and challenging to
98 specify in its entirety [7], two sources of uncertainty usually play a prominent
99 role in energy system modelling and beyond:

- 100 1. *Model Discrepancy (MD)*. Due to unavoidable modelling assumptions
101 and numerical approximations, the output simulated by a model f at
102 an input \mathbf{x} will be different from the value that would be observed in
103 the real world under the same physical conditions represented by \mathbf{x} .
104 This difference is referred to as model discrepancy.
- 105 2. *Observational Error (OE)*. It is a source of uncertainty intrinsic to all
106 physical measurements, with a magnitude that depends on the em-
107 ployed measurement device

108 Accounting for these two uncertainty sources is crucial to assess whether
109 a model has been successfully calibrated against observed data, and to make
110 sure that robust inference can be drawn from the calibrated model predic-
111 tions.

112 *2.2. Model Discrepancy in Building Energy Systems*

113 This section highlights sources of model discrepancy typical of building
114 energy models, arising from a wide range of assumptions that analysts need
115 to make

116 *2.2.1. Building Thermal Properties*

117 In a study of two office buildings in Australia [17], cooling set-points,
118 lighting power density, ICT infrastructure and its schedule were found to
119 impact the predictions of energy models the most. In other words, model
120 predictions were found to be particularly sensitive to small changes in these
121 parameters. Secondly, a major source of model discrepancy concerns the fab-
122 ric thermophysical properties, the fabric composition (i.e., non-homogeneity,
123 moisture content affecting thermal conductivity, etc.), as well as surface prop-
124 erties (radiative or convective characteristics). Most notably, the hypothesis
125 of mono-dimensionality of heat flow - which is fundamental to thermal resis-
126 tance (U-value) calculations in ISO 6946:2007 (Mohammed, can you please
127 provide the bibtex references for SO 6946:2007 and ISO 9869-1:2014, corre-
128 sponding to references number 5 and 8 of the word document?, ISO 9869-
129 1:2014, CIBSE and ASHRAE [5-8] (same for CIBSE an ASHRAE, references
130 number 6 and 7. I can't uniquely locate them on google scholar nor else-
131 where.) (Maybe check, whether the ASHRAE reference - number 7 - can be
132 replaced by my current number 3?) – remains a major simplification. In solid
133 bodies, heat travels in a diffused and 3-dimensional manner which no building
134 energy modelling platform replicates. Thirdly, for existing structures, notable
135 differences in thermal property measurements exist across seemingly uniform
136 building envelops [18–21]. These divergences are then compounded further
137 given that they are used in U-value calculations (either semi-stationary or
138 dynamic) that each produce a slightly different outcome [19] (disagreements
139 of up to 393% are reported between measured and calculated figures [18]).
140 While traditional solid masonry walls [22, 23] and floors [24] have been found
141 to perform better than model calculations, actual composite walls are re-
142 ported to perform worse than models suggest. A. Marshall et al. found
143 CIBSE model U-Values for brick walls, ceilings and doors to over-predict
144 performance by 30.3%, 15.5% and 9.9% respectively while calculated win-
145 dow properties showed much better matches with measured figures [25]. A
146 study of 57 properties found that while variations between similar wall types
147 and even within the same dwelling's walls existed, 44% of walls performed
148 better than CIBSE calculations, 42% were within acceptable bounds and
149 only 14% of sampled walls performed worse than calculations suggest [26],
150 however measured and calculated floor U-values were found to be in good
151 agreement. In summary difficulties in identifying the composite material and
152 their density and moisture content, internal and external air velocities, fabric

153 non-homogeneity and cumbersome nature of in-situ fabric studies are notable
154 reasons behind fabric U-value uncertainties.

155 *2.2.2. Weather and Occupant Activity*

156 The micro-climate around a building is by its nature random, with a
157 particular set of conditions never repeating itself over the building's lifetime.
158 It is affected by wind speed and direction, solar irradiance both outside and
159 within the building, external temperature and humidity, cloud cover and
160 ground albedo, atmospheric pressure. These all interact in close-to-chaotic
161 manner to dictate a micro-climate, which is then represented in a simplified
162 annual weather file for most energy simulations. Weather files are normally
163 extracted from decades of actual data (e.g. using Finkelstein-Schafer method)
164 but remain a major source of energy simulation uncertainty for the following
165 reasons: i) building energy performance is affected by future rather than
166 past weather; ii) a single weather file can hardly represent all meteorological
167 conditions; iii) most weather files are not measured at the site but at a
168 nearby meteorological station (often exposed airports), leading to inadequate
169 representation of urban heat island and sheltering effects [27].

170 Second to micro-climate, the occupant behaviour is drastically simplified
171 to inform the current generation of simulation tools. Temporal schedules usu-
172 ally describe the variations in occupant-related activities in a homogenous
173 and deterministic form. More accurate representations of the occupant's
174 dynamic interaction with the building remain an active area of research.
175 Advances in pervasive sensing and data collection have enabled better un-
176 derstanding of occupant presence and movements. Insights from improved
177 monitoring have informed two categories of occupant representation: implicit
178 and explicit models that predict stochastic occupant interventions with win-
179 dows, ICT and HVAC. Improvements are gained by seeking to understand
180 the underlying logic behind occupant interaction with its surrounding at a
181 behavioural level (explicit method) [28]. Challenges remain in adequately de-
182 tailed monitoring of occupant behaviour, the integration of future occupant
183 behaviour into simulation models and validation of the approach against lon-
184 gitudinal field studies. W. Tian et al. [29] note the need for further research
185 to provide clear guidance not only to characterise stochasticity but also on
186 the number of samples required to arrive at acceptable statistical power in
187 building energy analysis [29]. This remains a priority particularly if post
188 pandemic workplace attendance continues to evolve under a hybrid working
189 regime.

190 2.3. Limitations of Current Calibration Approaches in Building Energy Mod-
191 els

192 None of the approximations and assumptions discussed in Section 2.2
193 undermines the validity of a model. However, awareness of similar assump-
194 tions and approximations is key to quantify the discrepancies induced when
195 comparing model outputs and real-world observations.

In the context of building energy models, ASHRAE guidelines [3] are often followed to assess whether a model has been successfully calibrated against observed data [4, 5]. The guidelines make use of the following discrepancy measures between a sequence of N simulated model outputs S_i , and a corresponding sequence of observations M_i :

$$\text{MBE} = \frac{\sum(M_i - S_i)}{\sum M_i} \quad (1)$$

$$\text{CV(RMSE)} = \frac{\sqrt{\frac{1}{N} \sum(M_i - S_i)^2}}{\frac{1}{N} \sum M_i} \quad (2)$$

196 The model is considered calibrated if the relevant condition of the following
197 two is met:

- 198 (a) Hourly measurements: $-10\% \leq \text{MBE} \leq 10\%$ and $\text{CV(RMSE)} \leq 30\%$.
199 (b) Monthly measurements: $-5\% \leq \text{MBE} \leq 5\%$ and $\text{CV(RMSE)} \leq 15\%$.

200 While the above criteria are easy to check, their use to assess model
201 calibration presents some limitations. Firstly, acceptance thresholds are in-
202 dependent of the level of accuracy to which measurements are available, and
203 of the level of discrepancy between the model and reality. Moreover, formulas
204 (1) and (2) are not applicable to some of the physical quantities for which
205 model outputs and observations may be available.

206 Indeed, the two formulas measure the magnitude of the (average) bias
207 and root mean square error of the sequence $M_i - S_i$, relative to the (average)
208 measured value. This is meaningful for any quantity intrinsically defined on
209 a positive scale, such as energy consumption, but becomes less meaningful
210 for other quantities, such as temperature. Here, using different units (*e.g.*,
211 degrees Celsius or Kelvin) leads to seemingly different results. In reality, only
212 Kelvin should be used, but in this case criteria (a) and (b) would be easily
213 fulfilled due to the large denominator in (1) and (2).

214 In Section 3, we propose a statistical framework which overcomes these
215 limitations. The framework makes it easy to account for recognised un-
216 certainties when model and data are compared and it is applicable to any
217 quantity of interest (and, in fact, to several quantities simultaneously). More-
218 over, thanks to the use of emulators, the framework can explore the model
219 response at millions/billions of inputs, in a fraction of the time required for
220 physics-based models to generate output.

221 **3. Methodology**

222 This section presents statistical methods to quantify and include differ-
223 ent sources of uncertainty during model calibration. The overall process
224 we illustrate is referred to as History Matching (HM). The name was first
225 used within a study on hydrocarbon reservoirs [13] and indicates the attempt
226 to “replicate history” i.e., to produce model outputs that match historical
227 data. Since then, its statistical principles and methods have been success-
228 fully applied to a variety of fields (climate [15], biology [30], epidemiology
229 [12], cosmology [14] and more). However, in energy systems there has been
230 little to no application that are reported in the literature.

231 In this section we discuss the tools and statistical principles underlying
232 HM, illustrate the procedures and provide statistical packages to implement
233 it. We hope that our exposition, and the following illustration to a building
234 energy example study, will promote its uptake within the energy commu-
235 nity, and enable a robust treatment of uncertainty in calibration of computer
236 models.

237 The HM procedure sequentially rules out regions of the model’s input
238 space where, in light of all quantified uncertainties, model outputs cannot
239 match observed data. Emulators and implausibility measures are the two
240 key tools used to accomplish the task as follows:

- 241 • Emulators allow to predict the model output at choices of input param-
242 eters where the model has never been run, quantifying the uncertainty
243 of the prediction. Given their speed, emulators can be used to explore
244 the model parameter space thoroughly, even in higher dimensions.
- 245 • Implausibility measures quantify the distance between model outputs
246 and observations, in light of different sources of uncertainty.

247 Section 3.1 provides an overview of emulation, in simple but comprehen-
248 sive terms. Sections 3.2–3.4 then formalise the concepts of model discrepancy

249 and measurement error, introduce implausibility measures, and illustrate the
250 overall HM procedure, respectively.

251 *3.1. Bayes Linear Emulators as Fast Model Surrogates*

252 Evaluating the model’s response across its input space, and comparing
253 its outputs to observations, is a key part of calibration. However, running
254 a number of simulations that thoroughly cover the parameter space proves
255 unfeasible due to time and computational constraints. Emulators can be
256 employed to overcome this issue.

257 An emulator is a statistical surrogate of the model, predicting the model’s
258 outputs at inputs where the model has never been run. Emulators’ most ev-
259 ident advantages are speed and low computational requirements. A nominal
260 measure of this speed is that predictions of the output of interest at thou-
261 sands of different inputs can be obtained in one second on a personal laptop.
262 Moreover, the emulator prediction is accompanied by a precise uncertainty
263 statement of its accuracy, which lays the basis for model uncertainty analysis.

264 In this work we employ Bayes linear emulators (BLEs) to replicate the
265 behaviour of the example-study building energy model (see section 4). One
266 emulator is built for each output of interest. In this section, we give a brief
267 overview of the choices behind the construction of a BLE, referring the reader
268 to [6, 14] for further details. The interested reader is also referred to [31] for
269 a much more in-depth treatment of Bayes Linear principles.

270 We use the following form for $f(\mathbf{x})$, the output of the computer model at
271 input \mathbf{x} :

$$f(\mathbf{x}) = \sum_j \beta_j g_j(\mathbf{x}) + u(\mathbf{x}). \quad (3)$$

272 Equation (3) is made of two terms. The first one is a regression com-
273 ponent: it uses a linear combination of known functions $g_j(\mathbf{x})$ to model the
274 global behaviour of f across the input space. The second term, instead, is
275 meant to capture the local fluctuations of f , once the main regression com-
276 ponent has been accounted for.

277 Now suppose that the model has been run at a sequence of inputs $\mathbf{x}_1, \dots, \mathbf{x}_n$.
278 Each \mathbf{x}_i is called a design point and we will denote the corresponding output
279 by $y_i = f(\mathbf{x}_i)$. Then:

- 280 • The coefficient β_j in equation (3) can be estimated by fitting a linear
281 regression model with predictors $g_j(\mathbf{x})$ to the known pairs (\mathbf{x}_i, y_i) .

282 • Values of the residual process $u(\mathbf{x})$ will then be known at each design
 283 point: $u(\mathbf{x}_i) = u_i$. These values will approximately oscillate around
 284 0, with local patterns that the regression term has not been able to
 285 detect. We predict values of $u(\mathbf{x})$ via a BLE.

286 The idea to this last point is to model $u(\cdot)$ as a stochastic process, specifying
 287 a prior mean and covariance structure for it. In the Bayes linear framework,
 288 these are then adjusted to the observed values u_i at \mathbf{x}_i , yielding an adjusted
 289 mean prediction and uncertainty statement for $u(\mathbf{x})$, at any \mathbf{x} .

290 The prior specifications we make for u are as follows. We assume $u(\mathbf{x})$ to
 291 have mean zero at all inputs \mathbf{x} and we use a stationary kernel to model the
 292 covariance between the value of u at any two inputs \mathbf{x} and $\tilde{\mathbf{x}}$. A stationary
 293 kernel is one for which the above covariance only depends on the difference
 294 $\mathbf{x} - \tilde{\mathbf{x}}$. For common choices in emulation see Section 4.2.1 of [32]. In this
 295 work, we use the squared exponential kernel, leading to the following formula:

$$\text{Cov}[u(\mathbf{x}), u(\tilde{\mathbf{x}})] = \sigma_u^2 \exp\left(-\sum_k \left(\frac{x_k - \tilde{x}_k}{d_k}\right)^2\right) \quad (4)$$

296 The subscript k in x_k and \tilde{x}_k denotes the k^{th} component of the two inputs,
 297 *i.e.*, the k^{th} model parameter. The positive coefficient d_k (called correlation
 298 length, or length scale) measures the strength of correlation in outputs when
 299 the k^{th} parameter is modified. Finally, the coefficient σ_u^2 denotes the prior
 300 variance of $u(\mathbf{x})$ at all inputs \mathbf{x} .

301 Expression (4) has been written in terms of all components of the input
 302 vector \mathbf{x} . However, just some of the parameters often prove relevant to
 303 explain the great majority of the variability in u . We call these active parameters
 304 and, in practice, only use them in expression (4). The small variability
 305 left in u , due to the inactive parameters, is modelled as uncorrelated noise
 306 with constant variance ν^2 .

307 Notice that no distribution has been assumed for the process u . Within
 308 the Bayes linear framework, we only need to specify mean and covariance, as
 309 above. These mean and covariance are then adjusted to the values $u(\mathbf{x}_i) = u_i$,
 310 yielding an adjusted mean and variance for the value of $u(\mathbf{x})$ at any \mathbf{x} . Details
 311 of the adjusted mean and covariance formulas can be found for example in
 312 [6, 14, 31]. The above should however be sufficient to clarify the setting
 313 within which a BLE is built, and the decisions that need to be made. In
 314 particular, these are: the predictors g_j of the regression term; the variances
 315 σ_u^2 and ν^2 ; the active parameters; the correlation lengths d_k .

316 Although – put strictly - the above procedure builds a BLE for the “resid-
 317 ual” process u , once the regression prediction at \mathbf{x} is added, an overall pre-
 318 diction for $f(\mathbf{x})$ is obtained. Ultimately, this is the focus of interest. For this
 319 reason, in the following the term emulator will refer to the overall emulator
 320 of f . At any input \mathbf{x} , this provides a prediction $\hat{f}(\mathbf{x})$ of the unknown value
 321 $f(\mathbf{x})$, and a standard deviation $\hat{s}(\mathbf{x})$ quantifying the potential error of such
 322 prediction.

323 *3.2. Model Discrepancy and Observational Errors*

324 The concepts of model discrepancy (MD) and observational error (OE)
 325 are relevant to any uncertainty analysis linking computer simulations and
 326 measurements [13–16]. In this section we formalise them using a general
 327 framework. To relate the notation to the context of building energy models,
 328 the reader may think of: \mathbf{x} as a set of model parameters; $f(\mathbf{x})$ as the simu-
 329 lated consumption in the building under parameters \mathbf{x} ; y as the actual energy
 330 consumed in the building; z as the (imperfect) meter reading reporting the
 331 consumed energy.

332 Following[6], we assume that an appropriate choice of inputs exists, \mathbf{x}^* ,
 333 that accurately represents the values of the system’s parameters. We then
 334 link the simulator output $f(\mathbf{x}^*)$ at this parameters configuration and the
 335 real-world value y of the quantity being modelled, via the formula:

$$y = f(\mathbf{x}^*) + \varepsilon_{MD}, \quad (5)$$

336 where the model discrepancy term ε_{MD} accounts for the difference between
 337 the real and the simulated process. While the assumption is that just one
 338 quantity is being modelled, the same formulation can be considered if both
 339 y and $f(\mathbf{x})$ were vectors.

340 The unobservable real-system value y is usually estimated through a mea-
 341 surement z . We link these two quantities via the relationship:

$$y = z + e_{OE}, \quad (6)$$

342 where the term e_{OE} accounts for the observational error in the measurement.

343 The additive formulation in equations (5) and (6) is a simple but efficient
 344 way to model MD and OE, which also makes statistical inference tractable.
 345 Information about the quantities ε_{MD} and e_{OE} should in fact be sought in
 346 statistical form, rather than quantified as a single number.

347 Manufacturer guidelines are usually available to estimate the OE magnitude
348 (e.g., up to 5% of the measured value). Estimates of MD magnitude
349 may be more challenging to obtain. However, the modeller's knowledge of
350 the assumptions and approximations used within the simulator, alongside
351 literature research, usually provide guidance on the uncertainty effects of
352 different assumptions and can therefore lead to an overall estimate of MD.

353 In the example study discussed in this work, we consider two different
354 levels of MD (10% and 20% of the emulated value of $f(\mathbf{x})$), discussing how
355 the choice affects the results (Section 5). For a more detailed treatment
356 of MD, and the further distinction between internal and external MD, the
357 reader is referred to [6].

358 *3.3. Implausibility Measures*

359 Given a computer model f of a physical system and measurements z of the
360 system, this section outlines how to measure the likelihood that an input \mathbf{x}
361 to the model well represents the system's features. We do this by comparing
362 $f(\mathbf{x})$ and z . As we will see, acknowledging the involved uncertainties will
363 naturally lead to measure implausibility as opposed to likelihood.

364 The setting is as follows. A series of observations z_1, \dots, z_m of m different
365 quantities is available. Each quantity is modelled by the simulator: $f_j(\mathbf{x})$
366 denotes the simulator prediction of quantity j at input \mathbf{x} . For instance, in
367 the case of a building energy model, z_j and f_j may represent measured and
368 predicted consumption in month j . To make the notation lighter, in the fol-
369 lowing we initially drop the index j , reintroducing it later where information
370 from different measurements/model outputs will be aggregated.

371 The construction of implausibility measures is based on the following
372 reasoning. A choice of input parameters \mathbf{x} well represents the hidden system's
373 parameters (called \mathbf{x}^* in Section 3.2) if the following two quantities coincide:

- 374 i) The value assumed by the real system under conditions \mathbf{x} (eg, energy
375 consumed under building parameters as in \mathbf{x})
- 376 ii) The actual value y of the real system (eg, actual energy consumed in
377 the building).

378 None of the two quantities is available for general \mathbf{x} , but both can be esti-
379 mated, respectively by $f(\mathbf{x})$ (the model predicted consumption at \mathbf{x}) and z
380 (the measurement of y , such as a meter reading).

381 Hence, by comparing the difference $f(\mathbf{x}) - z$ to the tolerance allowed by
382 MD and OE, we can establish whether the quantities in i) and ii) may or

383 may not be the same. This leads us to consider the following quantity:

$$\frac{|f(\mathbf{x}) - z|}{\sqrt{\text{Var}(\varepsilon_{MD}) + \text{Var}(e_{OE})}} \quad (7)$$

384 High values of the the quantity above make the input \mathbf{x} *implausible* to match
385 the unknown system features \mathbf{x}^* having led to observe z .

386 Notice that expression (7) requires knowledge of the simulated value $f(\mathbf{x})$.
387 In practice, this is rarely readily available. However, if an emulator of f
388 has been trained, the emulator prediction $\hat{f}(\mathbf{x})$ can be used as estimate of
389 $f(\mathbf{x})$. To account for the substitution, the emulator variance $\hat{s}(\mathbf{x})^2$ should
390 be included in the denominator of (7). Therefore, reintroducing the index j
391 to indicate the particular output in question, we define:

$$I_j(\mathbf{x}) = \frac{|f_j(\mathbf{x}) - z_j|}{\sqrt{\text{Var}(\varepsilon_{MD}^j) + \text{Var}(e_{OE}^j) + \hat{s}_j(\mathbf{x})^2}} \quad (8)$$

392 the implausibility measure of input \mathbf{x} , with respect to quantity j . Finally, as
393 overall measure of implausibility with respect to all the m quantities under
394 consideration, we define:

$$I(\mathbf{x}) = \max_j I_j(\mathbf{x}). \quad (9)$$

395 The rationale behind the above definition is as follows. A high value of
396 $I(\mathbf{x})$ implies that at least one of the $I_j(\mathbf{x})$ is high. In other words, the input
397 \mathbf{x} has led to an output $f_j(\mathbf{x})$ that is highly incompatible with the observed
398 value z_j , even after accounting for the involved sources of uncertainty (MD,
399 OE, emulation). In this case, we say that the input \mathbf{x} is *implausible* to
400 lead to a match. Vice versa, we call \mathbf{x} *non-implausible*, since the involved
401 uncertainties make it possible for \mathbf{x} to lead to a match between model outputs
402 and observations.

403 The history matching algorithm, described in Section 3.4, sequentially
404 rules out all implausible inputs, identifying a final region where model out-
405 puts and observations can match in light of quantified uncertainties.

406 3.4. History Matching: The Algorithm

407 History matching (HM) proceeds in waves. The procedure is illustrated
408 below and summarised thereafter in a schematic algorithm.

409 In the first wave, some or all of the model outputs for which observations
 410 are available are emulated. The implausibility measure (9) is computed across
 411 the space, and all inputs \mathbf{x} for which $I(\mathbf{x})$ exceeds a given threshold are
 412 discarded as implausible. The remaining region, comprising the currently
 413 “non-implausible” points, is usually referred to as the Not-Ruled-Out-Yet
 414 (NROY) region.

415 The process is repeated, in consecutive waves. Crucially, at any new wave,
 416 additional simulations are run within the current NROY region, before new
 417 emulators are trained and an additional fraction of the region is discarded
 418 as implausible. The procedure terminates when the NROY region does not
 419 shrink any further, or it has become empty (all inputs have been ruled out
 420 as implausible).

421 The algorithm below summarises the steps. Comments on the procedure
 422 and guidance on some of the choices follow in Section 3.5.

423 *History Matching Algorithm*

424 Choose a positive threshold T (guidance provided in Section 3.5). Define a
 425 sequence of NROY regions $\mathcal{R}^0 \supseteq \mathcal{R}^1 \supseteq \dots \mathcal{R}^k \supseteq \mathcal{R}^{k+1} \dots$ as follows:

- 426 1. Initial Step ($k = 0$)
 - 427 (a) Let \mathcal{R}^0 be the whole input space. Perform a sequence of model
428 runs in \mathcal{R}^0 , and build emulators of some of the quantities for which
429 observations are available.
 - 430 (b) Let I^0 be the implausibility measure in equation (9). Define the
431 first NROY region as $\mathcal{R}^1 = \{\mathbf{x} \in \mathcal{R}^0 : I^0(\mathbf{x}) < T\}$.
- 432 2. Iterative Step ($k \geq 1$)
 - 433 (a) Run additional simulations within \mathcal{R}^k . Decide which outputs to
434 emulate, and build emulators for them based on old and new runs
435 within \mathcal{R}^k .
 - 436 (b) Compute the implausibility $I^k(\mathbf{x})$ over \mathcal{R}^k , and identify \mathcal{R}^{k+1} as
437 $\mathcal{R}^{k+1} = \{\mathbf{x} \in \mathcal{R}^k : I^k(\mathbf{x}) < T\}$.
- 438 3. Stop when $\mathcal{R}^{k+1} = \mathcal{R}^k$ or \mathcal{R}^{k+1} is empty.

439 The flowchart in Fig. 1 illustrates the overall methodology.

440 *3.5. Comments on the Procedure and its Strengths in Handling Uncertainties*

441 We comment here on some of the choices needed to perform HM and on
 442 the overall strengths of the procedure, particularly in handling uncertainties.

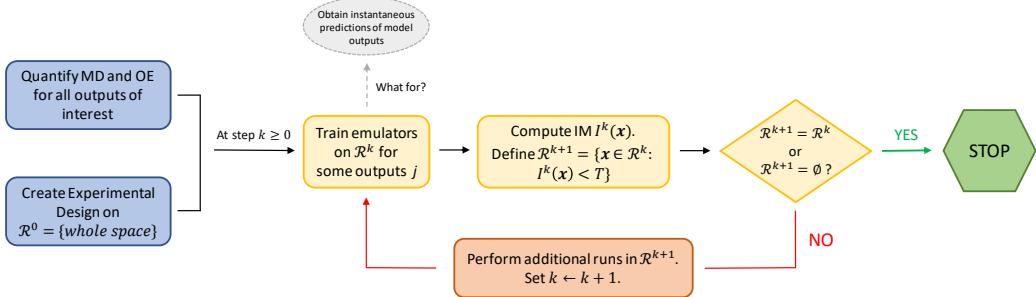


Figure 1: Flow chart illustrating the History Matching (HM) procedure. The initial steps (blue) consist in quantifying the main sources of uncertainty and designing a first set of runs over the whole space. The central part of the diagram (yellow) describes the steps of a typical wave, which aims to reduce the volume of the region comprising currently non-imausible inputs. Consecutive waves are repeated, till the region cannot be further reduced, or it is empty.

443 The choice of the threshold T reflects how large an implausibility we are
 444 prepared to accept, before ruling out an input as implausible. A common
 445 choice is $T = 3$, based on Pukelsheim's 3σ rule [33]. See [14] for more
 446 details. However, a larger threshold may be chosen if several quantities are
 447 being history-matched. Indeed, mainly due to emulation error, for each fixed
 448 output j there is some small probability that $I_j(\mathbf{x}) > T$ even if \mathbf{x} was the
 449 “best” input \mathbf{x}^* . This probability increases if we consider the event that *at*
 450 *least* one j leads to $I_j(\mathbf{x}) > T$, that is, the event $I(\mathbf{x}) > T$. Therefore, a
 451 larger threshold T ensures that this probability of error is kept low.

452 As highlighted by point (2a) of the algorithm, at each wave, only some
 453 of the model outputs are used to compute the implausibility measure. This
 454 is a key feature of HM. Especially in early waves, some of the outputs may
 455 be difficult to emulate over large parts of the space, typically because the
 456 model’s behaviour varies significantly across different input regions. As the
 457 search of non-imausible inputs is narrowed down to much smaller regions,
 458 outputs that were difficult to emulate may behave more uniformly within the
 459 region of interest. They can therefore be emulated precisely and included in
 460 the definition of $I(\mathbf{x})$ to rule out further implausible regions.

461 Also notice that, especially in earlier waves, not ruling out an input
 462 ($I(\mathbf{x}) \leq T$) does not suggest that the input leads to a match with the obser-
 463 vations. Indeed, the implausibility of an input \mathbf{x} may be low as a consequence

464 of large emulator uncertainty (term \hat{s}_j in equation (8)), rather than actual
465 proximity between the emulated prediction and the observation. However,
466 emulator uncertainty is typically reduced between waves, because the emu-
467 lators are trained on smaller regions (where the simulator may behave more
468 smoothly) and on more design points. Reducing emulator uncertainty leads
469 to larger implausibilities and, therefore, to more inputs being identified as
470 implausible.

471 This is, in fact, the rationale behind proceeding in multiple waves and
472 a key strength of HM. The decision of rejecting an output as implausible
473 is taken in light of all quantified uncertainties at the given stage, including
474 MD and OE. In this regard, notice that expression (8) can be easily adjusted
475 to include additional sources of uncertainty in its denominator, whenever
476 these are recognised and quantified in the problem in question. Moreover,
477 by only involving differences and measurements of variability, implausibility
478 measures can be computed for any quantity for which measurements are
479 available, even if these are not measured on a positive scale (*e.g.*, Celsius
480 degrees for temperature). This marks a crucial difference with respect to
481 formulas (1) and (2).

482 In most cases, one or two waves are enough to rule out as implausible
483 a great percentage of the space. In fact, as shown in the example study of
484 Section 4, one wave may even be enough to identify the final non-implausible
485 region, if the emulator uncertainties are negligible with respect to MD and
486 OE. More generally, however, multiple waves may be needed. In R, the hmer
487 package [34] can be used to implement emulation and HM. The package
488 guides the researcher into multiple HM waves, by automatically suggesting
489 a new design for each wave.

490 We conclude this section with a consideration. By designing model runs
491 only where needed, HM allows the researcher to focus sequentially on the
492 region of interest, while keeping the overall number of runs low. This is
493 crucial in medium and high dimensions. In these cases, the region where a
494 match is possible may only represent a tiny fraction of the original space.
495 A more naïve yet not uncommon approach, where a number of simulations
496 are run across the space and compatibility with observations is checked on
497 these, would almost surely miss the region, even if the number of runs was
498 very large. HM allows instead to identify the region precisely (even in high
499 dimensions) thanks to the sequential “refocussing” procedure.

500 The event where, at some wave, all inputs have been discarded as im-
501 plausible, is possible and indeed very informative. Causes for the lack of

502 mismatch between model and data may be multiple, with typical ones as
503 follows: i) original parameter ranges are incorrect, ii) there is a problem with
504 the data, iii) uncertainties are higher than estimated, iv) the model dynam-
505 ics has a flaw. It is of course the researcher's task to step back and analyse
506 what causes the mismatch, intervening on the model or in the assessment of
507 uncertainties if necessary.

508 4. Example Study

509 This section describes the example study building used to illustrate the
510 methodology outlined in Section 3. The building is a single, detached dwelling
511 for which a model has been created in EnergyPlus (E+). Longitudinal energy
512 consumption and temperature data is collected to calibrate the model.

513 4.1. Building and its Energy Usage

514 The detached two-story masonry construction used as case study was
515 built in 1994. Two occupants are the only residents of the dwelling and were
516 asked to archive their gas cooker and shower usage each day across an an-
517 nual cycle. Given a very predictable pattern of occupancy (both occupants
518 had 8am-5pm working commitments), it was possible to limit the stochastic
519 nature of occupant activity as far as practically manageable and use deter-
520 ministic schedules to represent occupant interventions with the building and
521 its energy system. The building (with a gross area of 168.66m² and 19.73m²
522 of unheated space) is located in a UK built-up urban surrounding and is only
523 partly shaded on its west elevation by another adjacent property (shading
524 represented in the model).

525 Across the monitoring year (2016) the property had an observed annual
526 gas (15,381 kWh) and electricity (2,991 kWh) consumptions that respec-
527 tively correspond to high and medium UK domestic consumption values [23]
528 (Mohammad: reference 23, can you please provide it in bibtex format? Same
529 with ref 24 and 25 below please. I cannot clearly identify them in google
530 scholar). Occupants utilised shower facilities at a measured flow rate of 4.37
531 l/min and recorded on average eight 20-minute showers per week correspond-
532 ing to an average of 50 l/person/day (occupants used cold water over wash
533 basin and dishwasher supplied with cold feed only). These recorded values are
534 below UK average domestic hot water usage (reported as 142 l/person/day
535 [24] and 122 l/person/day [25]), but primarily reflect occupants' heavy use

536 of gym washing facilities. Gas cookers (containing 3kW and 5kW hubs) were
 537 used on average 4 times a week for 1 hour per cooking session.

538 *4.2. The model: Co-dependency of Energy and Temperature Predictions*

EnergyPlus (E+) is a collection of dynamic modules each simulating different environmental, climatic and operational conditions that define either the flow or the stored quantity of energy within building internal zones. The core of the programme is a heat-balance equation that is solved for all zones using one of three methods (3rd order backward difference, Euler method or analytically) to converge zone loads and resultant temperatures to within a pre-defined tolerance (using a predictor/corrector process). The uniqueness of E+ lies in it being a physically-based modelling solution. It oversees i) a simultaneous calculation of radiative and convective heat and mass transfer processes, ii) adsorption and desorption of moisture in building elements and iii) iterative interactions of plant, building fabric and zone air. This integrated and simultaneous simulation process is completed via several modules (and overseen by E+ simulation manager), with understandably multiple first-principle-based equations that are solved simultaneously and/or iteratively. This makes it very difficult to bring a sharp focus on any single or sets of expressions where model prediction uncertainties lie. However, the zone air heat balance equation is one of the primary mechanisms that describes the connected nature of heat gains/losses within a zone, the corresponding plant duty to offset them and the resultant zone mean air temperatures:

$$C_z \frac{dT_z}{dt} = \sum_{i=1}^{N_{sl}} \dot{Q}_l + \sum_{i=1}^{N_{sur}} h_i A_i (T_{si} - T_z) + \sum_{i=1}^{N_{zones}} \dot{m}_l C_p (T_{zi} - T_z) + \dot{m}_{inf} C_p (T_\infty - T_z) + \dot{Q}_{sys}, \quad (10)$$

539 whereby $C_z dT_z/dt$ is the rate of change of thermal energy stored in the zone
 540 air, $\sum \dot{Q}_l$ is the sum of convective internal loads (**Mohammad, can you check**
 541 **that in the first and third sums, the summing index is meant to be lowercase**
 542 **L?**), $\sum h_i A_i (T_{si} - T_z)$ is convective heat transfer from the zone surfaces,
 543 $\sum \dot{m}_l C_p (T_{zi} - T_z)$ is the change of the room air enthalpy as a result of zone
 544 air mixing, $\dot{m}_{inf} C_p (T_\infty - T_z)$ is the infiltration heat transfer and finally \dot{Q}_{sys}
 545 is the HVAC system input to achieve target temperature. Essentially each
 546 item on the right-hand side of equation (10) indicates a change of enthalpy



Figure 2: Left panel: power monitor used to characterise household appliances. Right panel: AC current sensor, monitoring transmitters and temperature sensors deployed in the case-study building.

547 due to environmental perturbations, while the left-hand side describes how
 548 these perturbations impact the zone air temperature and its enthalpy.

549 E+ assumes uniformity in i) zone air and surface temperatures, ii) in-
 550 coming long/short-wave/diffused radiations and iii) surface properties (as
 551 opposed to treating radiation in a direct or point-based manner). It is reason-
 552 able to regard zone air temperature as the interconnection where conductive,
 553 radiative and convective heat balance and mass transfers are realised. This
 554 underpins our model validation approach in which energy and temperature
 555 data are essential.

556 *4.3. Data Collection*

557 A proprietary set of environmental and energy sensors were deployed to
 558 compile electricity and zone temperatures (Fig. 2). To reduce measurement
 559 uncertainty, each of the two target zones was equipped with two separate
 560 air temperature sensors at 1.3m above floor level and set to log data at 30s
 561 intervals to achieve a moderated average. Therefore, space temperature was
 562 recorded by 4 battery-powered sensors: two positioned in the south-facing
 563 master bedroom, and two in north-facing kitchen. Overall kitchen and master
 564 bedroom temperature sensors had total annual losses of 5.7% and 2.7% that
 565 required imputation. Each missing temperature cell was imputed by the
 566 average of the previous and successive available cells.

Table 1: Input parameter variations explored in the batch runs used to train our Bayesian emulators. Mohammad, many of these values need a careful checking between the two of us, see email :)

	Heating Setpoint (°C)	Boiler Seasonal Efficiency (%)	Ext. wall U-value (Unit)	Roof U-val (Unit)	Floor U-val (Unit)	Infiltration rate (ach)	DHW Consumpt. (l/pers/day)	Cooking (%)
Short name	V1	V2	V3	V4	V5	V6	V7	V8
Base-model Value	17.5	65	0.544 (out of range)	0.213	0.337 (out of range)	0.20	to be discussed	missing
Range	[17.5, 20.5]	[60, 75]	[0.04, 0.063]	[0.15, 0.21]	[0.045, 0.055]	[0.2, 0.95]	$[6.15 \times 10^{-6}, 2.2 \times 10^{-5}]$	[1.05, 6.3]

567 Gas consumption data was manually recorded at monthly basis using
 568 mains gas meter. Electricity consumption was logged at 10s intervals using
 569 two mains-powered clip-on current sensors on the incoming live cable pro-
 570 viding two identical readings. The average of the two (in the end identical)
 571 readings formed the measured power usage. Gas and electricity data required
 572 no imputation.

573 In order to parameterise the energy model more accurately, a plugin power
 574 monitor was also used to characterise instantaneous and time-averaged con-
 575 sumption of the main electrical devices (TV, washing machine, ICT).

576 4.4. Model Input, Output and Range

577 By consulting manufacturers specification and the house builder's litera-
 578 ture, a detailed set of parameter inputs for the model were compiled. Lower
 579 and upper bands were derived from scientific literature and used to dictate
 580 the size of associated variations explored in batch-runs (Table 1). Variations
 581 in internal and external air velocities are less pronounced in ground floors,
 582 but more significant for walls and roofs. Therefore, smaller floor uncertainty
 583 margin are reported in literature (Tables 1 and 2: Mohammad pls check
 584 correct referencing).

585 Manufacture's literature for glazing (installed in 2009) stated respective
 586 G and U-values of 0.69 and 1.79 W/m²K. Respective error bands of ±5%

Table 2: Parameter inputs for energy model development of the case-study building.
(Dario-Mohammad: discuss caption and other red points in zoom call)

Parameter	Description
Heating	Natural gas boiler serving a radiator central heating system
Heating setpoint (setback)	19°C (16°C)
Heating schedule	02–11, 16–24
Ventilation	Natural ventilation (mechanical extract to family bathroom and en suite)
Ventilation rate	Highly stochastic, controlled by via openable windows
Gas boiler seasonal efficiency	65% (15 year old non-condensing gas-fired system boiler – 77°C/55°C F+R)
DHW consumption	0.59 litre/m ² /day
Cooling setpoint (setback)	Uncontrolled
Nominal lighting power density	1.4 W/m ² (manually controlled) to achieve 200 lux
Number of occupants	2 in total
Internal gains ¹	6 W/m ²
Gross (conditioned) area	168.66m ² (148.93m ²)
Observed annual gas (electricity) consumption (2016)	15,381 kWh (2,991 kWh)
There were more, but are already present in Table 1	

587 and $\pm 2\%$ for G and U- values altered the gas consumption by ± 2.05 kWh
588 ($\pm 0.013\%$). Given its negligible nature, the error bands of the glazing were
589 discounted in batch simulations. Given opaque fabric U-value impact on
590 energy consumption [35], uncertainty bands are derived from literature and
591 imposed on opaque U-values (Table 2).

592 Estimation of actual building infiltration rates requires convoluted air
593 permeability tests. Table 4.16 of CIBSE guide A [6] (For ref 6 see one of above
594 notes) outlines a range of 0.25 to 0.95 air change per hour (ACH) for various
595 2-story buildings below 500m^2 with a value of 0.5 ACH describing typical
596 constructions similar to the case-study building. Therefore 0.5 ACH with a
597 range of 0.25-0.95 ACH are used (meaning of 0.5? discuss with Mohammad).
598 Local weather files compiled by a weather station approximately 3 miles away
599 from the site was used to support the model development [27] (M: can you
600 pls provide bibtex version of old reference 27).

601 We therefore consider eight uncertain parameters (Table 1), and use
602 monthly gas and electricity consumption and hourly temperature data to
603 history-match them. In the following, we will use the parameter to refer to
604 any of the eight uncertain quantities in Table 1, and the term input to refer
605 to a specific configuration of the eight parameters $\boldsymbol{x} = (x_1, \dots, x_8)$.

606 5. Results

607 This section discusses the results of the application of the statistical prin-
608 ciples and methods of Section 3, to the example study of Section 4. We build
609 emulators of the model outputs of interest and history-match these outputs
610 to identify a region of non-imausible inputs.

611 5.1. Experimental Design and Simulated Results

612 To build emulators of the model outputs, results of a small number of
613 model runs are needed. An informal rule of thumb suggests using about 10
614 times as many runs as the number of varied parameters (eight in our case)
615 [36]. In our case the simulator was moderately fast to evaluate, and a total of
616 $n = 1,000$ simulations were run (taking about 16 hrs in E+). Note that this
617 is substantially more than what is typically needed for eight inputs. While
618 the results presented here are obtained by exploiting the information from
619 all 1,000 runs, we have also simulated a more typical case where only 100
620 runs (randomly selected among the original 1,000) were used. This however
621 yields no difference to the final results (see Section 6).

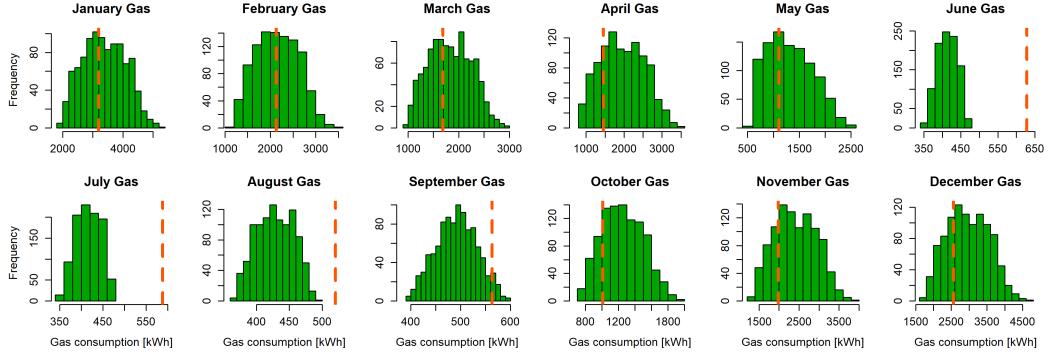


Figure 3: Distribution of simulated monthly gas consumption, for the experimental design used in this study ($n = 1000$ inputs). The dashed vertical line in a plot denotes the observed consumption for that month. Note that, for Jun-Jul-Aug, this lies outside of (and often far from) the range of simulated consumptions.

622 The experimental design (the set of n inputs for which the simulator is
 623 run) was chosen through Latin Hypercube Sampling (LHS). In two dimen-
 624 sions, the technique selects n points in a square so that exactly one point
 625 falls in each row and in each column of an $n \times n$ uniform grid of the square.
 626 The same idea governs higher dimensions. For this reason, LHS is commonly
 627 used in the design of computer experiment analyses [14, 37, 38], to identify
 628 points in a hypercube that fill the space well and are not too close to each
 629 other [39]. LHS is easily performed in R via the “lhs” package.

630 5.1.1. Gas

631 Fig. 3 shows simulated and observed values of monthly gas consumption
 632 for the 1,000 design runs. A distinction between summer and non-summer
 633 months can be noticed: in June, July and August, all model runs considerably
 634 underestimate the observed consumption. This, in principle, does not rule
 635 out the possibility that other inputs within the explored space (an eight-
 636 dimensional hypercube with side ranges as in Table 1) may lead to a match.
 637 In this case, however, especially in June and July, it is evident that the range
 638 of simulated outputs is small with respect to the distance between these and
 639 the observed consumption: a simple linear model confirms indeed that a
 640 match within the explored space cannot be achieved for these months.

641 As illustration of the proposed methodology, in the following we will look
 642 for inputs that match all nine non-summer monthly gas consumption si-
 643 multaneously, excluding the three summer conditions from the match. It is

644 important to note that in a real case study, the model's underestimation of
645 summer consumptions should be identified and acted upon, prior to calibra-
646 tion/history matching. While the most likely reason for summer months
647 deviation is the inability of deterministic schedules to capture stochastic
648 aspects of energy governance in the household (*i.e.*, additional use of hot
649 water/cooking in these months), our intervention with the model to reflect
650 historical occupant behaviour for which we had little certainty would have
651 been speculative. Greater insight into seasonal energy use variations would
652 have allowed one or more parameters to be added as additional input vari-
653 ables to be history-matched.

654 *5.1.2. Electricity*

655 As described in Section 4.4, the model also predicts monthly electricity
656 consumption. Within the explored ranges of the parameters, however, all
657 months display almost no variability in output, with an observed consump-
658 tion far outside the range of simulated values. For the three summer months
659 (Jun–Aug), no variation at all is shown in simulated consumption among the
660 n runs. This is partly an outcome of the deterministic electrical load speci-
661 fication. More importantly though, this is due to the fact that any changes
662 with the 6 variables in Table 1 except for heating and DHW (V1 and V7)
663 will only affect gas consumption. This makes the variability of electricity
664 consumption very limited and bound to a small deterministic input schedule
665 (e.g., a fixed 6 W/m² for internal gains as dictated by field measurements).

666 Given the illustrative role of the example study within this work, we will
667 not attempt to match the electricity consumptions in the methodological
668 illustration that follows. We will instead consider the nine conditions coming
669 from the observed monthly gas consumption in non-summer months, to which
670 we will add temperature constraints in Section 5.5.

671 *5.2. Emulation of Monthly Gas Consumption*

672 For each of the nine non-summer months, the $n = 1,000$ simulations pro-
673 vide a dataset of n pairs (\mathbf{x}_i, y_i) where \mathbf{x}_i represents one configuration of the
674 8 input parameters, and y_i is the simulated gas consumption associated with
675 it. We use the dataset to build an emulator of that month's gas consumption.
676 To that aim, the eight input parameters are linearly rescaled so that each of
677 them ranges in the interval $[-1, 1]$.

678 From Section 3.1, recall that several choices (on covariates, correlation
679 lengths, prior variances) have to be made when building an emulator. To

680 validate these choices, we split the dataset as follows:

- 681 • Training set (700 runs): used to train the emulators.
- 682 • Evaluation set (150 runs): used to decide on the values of the emulator
683 hyperparameters, by comparing the emulator’s performance on this set
684 to the known simulator’s outputs.
- 685 • Test set (150 runs): used to test the previously built emulators on a
686 completely new set of runs not used for training and evaluation.

687 Following the notation of Section 3.1, the prediction of the simulated con-
688 sumption at an input \boldsymbol{x} is computed as sum of: i) a linear regression part, ii)
689 a prediction of the regression residual at the point \boldsymbol{x} , in the form of a Bayes
690 linear emulator. Details on each of the two parts follow.

691 *5.2.1. Linear Regression*

692 The only choice to be made in building a linear regression model con-
693 cerns the predictors to use, *i.e.*, the functions g_i in equation (3). For all the
694 months of interest, a preliminary exploration reveals that the response y (gas
695 consumption) is very well explained as quadratic function of the eight input
696 parameters, which in the following we call V_1, \dots, V_8 for convenience. Thus,
697 we proceed as follows.

698 Let \mathcal{P} be the set of all mutually orthogonal linear, quadratic and inter-
699 action terms of V_1, \dots, V_8 .² For a given integer k , consider the linear model
700 with highest adjusted coefficient of determination (adj. R^2) and consider the
701 smallest integer k for which the above linear model with k predictors has an
702 adj. R^2 exceeding the threshold.

703 In our case, 10 predictors yield models with notably high adj. R^2 , about
704 0.999 across all months. To keep the approach uniform among months, we
705 consider the linear model with the “best” 10 predictors for each month (Ta-
706 ble 3). As pointed out above, however, a different number of predictors for
707 each output may be considered on other case studies.

708 Note that, in general, such a high R^2 should raise concerns of overfit-
709 ting. In our case, the concern can be ruled out by observing that only 10

²Obtained in R via the command `poly(X, deg=2)`, where X is the 700×8 matrix whose rows are the training inputs. For 8 parameters, there is a total 44 linear, quadratic and interaction terms.

Table 3: Properties of the gas consumption emulators. For each month, from left to right: covariates used to build the linear regression model; adjusted R^2 of the linear model; variance of the residuals; active parameters used in the covariance function; value of the correlation lengths (same for all active parameters). In the predictor's column, the * symbol denotes the combination of all linear and interaction terms: $a * b * c = \{a, b, c, ab, ac, bc\}$.

	Predictors	Adj. R^2	σ^2	Act. Params	d
Jan	$V_1 * V_2 * V_6, V_3, V_4, V_2^2, V_6^2$	0.9998	85.94	V_1, V_2, V_3, V_6, V_8	0.8
Feb	$V_1 * V_2 * V_6, V_3, V_4, V_2^2, V_6^2$	0.9997	52.30	V_1, V_2, V_3, V_7, V_8	0.65
Mar	$V_1 * V_2 * V_6, V_3, V_4, V_2^2, V_6^2$	0.9997	50.14	$V_1, V_2, V_3, V_6, V_7, V_8$	1.3
Apr	$V_1 * V_2 * V_6, V_3, V_4, V_2^2, V_6^2$	0.9998	50.04	V_1, V_2, V_3, V_6, V_8	1
May	$V_1 * V_2 * V_6, V_3, V_8, V_1^2, V_6^2$	0.9989	206.45	$V_1, V_2, V_3, V_4, V_6, V_8$	1
Sep	$V_1 * V_2 * V_6, V_3, V_4, V_8, V_6^2$	0.9994	0.91	$V_1, V_2, V_3, V_6, V_7, V_8$	1.2
Oct	$V_1 * V_2 * V_6, V_3, V_4, V_8, V_6^2$	0.9996	24.40	$V_1, V_2, V_3, V_6, V_7, V_8$	1.4
Nov	$V_1 * V_2 * V_6, V_3, V_4, V_2^2, V_6^2$	0.9998	57.35	$V_1, V_2, V_3, V_6, V_7, V_8$	1.2
Dec	$V_1 * V_2 * V_6, V_3, V_4, V_2^2, V_6^2$	0.9998	68.77	$V_1, V_2, V_3, V_6, V_7, V_8$	1.3

710 predictors have been used to explain 700 observations. The high coefficient
711 of determination mirrors an intrinsically quadratic model dynamics.

712 *5.2.2. Emulators of the Residuals*

713 To build an emulator of the regression residuals of each month's consump-
714 tion, choices about the following quantities are to be made: active parameters
715 \mathbf{x}_A , correlation lengths d_i , prior emulator variance σ_u^2 an nugget variance ν^2
716 - see equation (4) and text thereafter. We proceed as follows.

- 717 • **Active parameters \mathbf{x}_A :** To identify them, we look at the most signif-
718 icant (in the sense discussed in Section 5.2.1) second- and third- order
719 terms in a linear regression model of the residuals. Variables appear-
720 ing by themselves with a high t -value ($t > 8$) are included as active
721 parameters. The inclusion of variables appearing alone with a lower
722 t -value or in interaction with other variables is instead considered case
723 by case, according to the emulator performance on the evaluation set
724 – see Section 5.2.3.
- 725 • **Correlation lengths d_i :** Once the active parameters are chosen, the
726 same correlation length d is used for all of them. The value of d at
727 each month is chosen by assessing the emulator's performance on the
728 evaluation set and is reported in Table 3. Notice that, to attain a
729 similar level of correlation across the space, higher correlation lengths
730 are used when a higher number of active parameters is present.
- 731 • **Prior variances σ_u^2 and ν^2 :** Let σ^2 denote the variance of the regres-
732 sion residuals being fitted. Hence, we set $\sigma_u^2 = 0.95\sigma^2$ and $\nu^2 = 0.05\sigma^2$.

733 Table 3 provides details of all choices made to build each of the nine
734 emulators, including the ones concerning the regression line. We discuss
735 validation of the emulators in the next section.

736 *5.2.3. Emulator Validation and Performance*

737 The active parameters \mathbf{x}_A and correlation lengths d in Table 3 are chosen
738 based on the emulator's performance on the evaluation set. The latter con-
739 sists of 150 pairs (\mathbf{x}_i, y_i) not used in the emulator's training, where each y_i
740 is the simulated output at input \mathbf{x}_i . At each such \mathbf{x}_i , the emulator provides
741 a prediction \hat{y}_i of the simulated output and an uncertainty statement about

742 the prediction in the form of a standard deviation $\hat{\sigma}_i$. One way to assess the
 743 emulator's performance at \mathbf{x}_i is thus to define:

$$\varepsilon_i = \frac{\hat{y}_i - y_i}{\hat{\sigma}_i}. \quad (11)$$

744 We call ε_i the emulator standardised error at \mathbf{x}_i . It represents the number of
 745 standard deviations which separate the emulator prediction from the known
 746 simulated output.

747 As we make no distributional assumption on the emulator, we can appeal
 748 to Pukelsheim's 3σ rule [33] to constrain expected values of ε_i : the result
 749 states that, for any continuous unimodal distribution, at least 95% of the
 750 probability mass lies within 3 standard deviations from the mean. Therefore,
 751 by only assuming unimodality of the emulator distribution, we should expect
 752 at least 95% of the ε_i to lie between -3 and 3.

753 For each of the nine emulators, we consider the plot of ε_i versus the
 754 predictions \hat{y}_i . Given a good emulator, such a plot should be characterised
 755 by an approximately random scattering of the points around the line $\varepsilon = 0$,
 756 with about 95% of them in modulus less than 3 due to Pukelsheim's rule. We
 757 assess visually both these properties for different choices of active parameters
 758 and correlation lengths, choosing the ones which return plots with the desired
 759 properties.

760 We tend to be slightly conservative in this phase, by choosing correlation
 761 lengths which generally yield more than 95% of $|\varepsilon_i|$ less than 3 . This is
 762 to prevent making choices tailored to the specific points used for validation.
 763 Once the parameters \mathbf{x}_A and d are chosen, we compute the standardised
 764 errors in equation (11) on the 150 elements of the test set, to check that no
 765 anomaly shows up on a set of points never used in training or evaluation.

766 Finally, we note that the emulators of the nine monthly gas consumptions
 767 are remarkably precise in their predictions. This is a consequence of both an
 768 excellent regression fit to the data, and of a further reduction in the residual
 769 uncertainty due to the emulators. Quantification of this reduction is provided
 770 in Table 4, which shows the original variance of the residuals and an empirical
 771 95% confidence interval (CI) of the emulator variance on the 150 test points.
 772 The CIs show that the uncertainty in the emulator predictions is very small
 773 compared to the range of simulated outputs (Fig. 4).

Table 4: For each month: variance of the regression residuals to which the emulator is fitted (first row); empirical 95% confidence interval of the emulator variance on the 150 test points (second and third row).

	Jan	Feb	Mar	Apr	May	Sep	Oct	Nov	Dec
σ^2	85.9	52.3	50.1	50.0	206.4	0.91	24.4	57.4	68.8
[2.5%,	6.6	5.9	3.1	3.2	16.8	0.06	1.5	3.8	4.3
[97.5%]	31.0	28.0	9.3	10.0	68.9	0.20	3.8	12.8	12.7

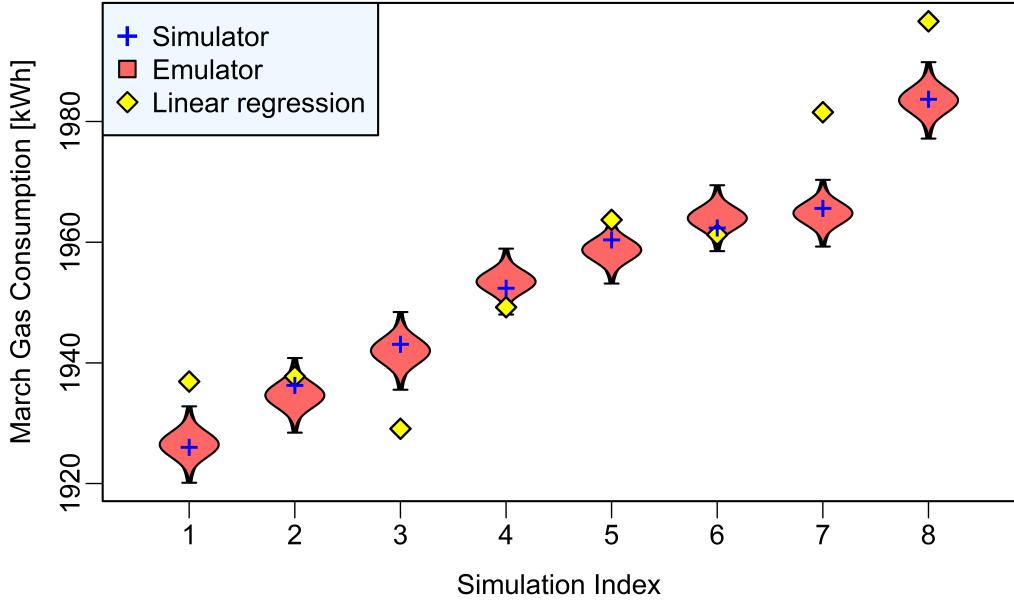


Figure 4: Comparison between point regression predictions and emulator predictions with associated uncertainty, for a sequence of randomly selected eight consecutive simulated March gas consumptions of the test set. The bell-shaped curves cover a range of 3 emulator standard deviations up and down the emulator prediction; the simulated value (+) always lies within this range, while the regression prediction may lie substantially further. Results are displayed for eight simulations to have a small range in the y -axis and allow to appreciate detail.

774 5.3. *Observational Error and Model Discrepancy*

775 The validated emulators can be used as probabilistic surrogate of the
776 original model, to predict the monthly gas consumption at any input \mathbf{x} .
777 We use the implausibility measure described in Section 3.3 to compare the
778 prediction with the observed consumption z . The observational error (e_{OE})
779 and model discrepancy (ε_{MD}) are set as follows:

- 780 1. The OE is set to 5% of the observed value z , in agreement with the
781 largest manufacturer's accuracy band ($\pm 5\%$ for power monitor).
- 782 2. For illustrative purposes, we set MD to either 10% or 20% of the emu-
783 lated consumption and compare results in the two cases.

784 Note that, as discussed in Section 2.2, accurate estimation of model discrep-
785 arity requires careful statistical analysis and possibly additional model runs.
786 However, an order of magnitude between 10 and 20% of the simulated val-
787 ues is likely to represent a good estimate in many energy applications. If
788 there are reasons to believe that model discrepancy is much higher, then the
789 possibility of revisiting the model itself should be considered, by possibly
790 including in the model key factors affecting the simulated dynamics.

791 5.4. *History Matching*

792 The HM procedure has been described in Section 3.4. In our case, at an
793 input configuration \mathbf{x} :

- 794 1. We first compute the implausibility measures $I_M(\mathbf{x})$ for months M, as
795 in equation (8);
- 796 2. We then define the overall implausibility $I(\mathbf{x})$ as in equation (9), $I(\mathbf{x}) =$
797 $\max_M I_M(\mathbf{x})$;
- 798 3. Hence, we consider an input as non-implausible if $I(\mathbf{x}) < 4$.

799 As explained in Section 3.5, choosing a higher threshold than the more
800 common $T = 3$ ensures that the probability of incorrectly rejecting an in-
801 put \mathbf{x} as implausible is kept small, whenever several observations are being
802 matched simultaneously (nine in our case). Also notice that an alternative
803 approach may be to define $I(\mathbf{x})$ as the second- or third-highest value of $I_M(\mathbf{x})$
804 over all months, while keeping a relatively low threshold T . However, if a
805 particular month proves more challenging to match, this approach risks to
806 classify several points as non-implausible, although those points are highly
807 unlikely to match the observed consumption for that month.

808 In our example study, the percentage of space classified as non-imausible
809 after one wave is 0.30% when 10% MD is used, and 19.52% for 20% MD.
810 These percentages have been estimated by computing the proportion of non-
811 implausible points on a sample of size $N = 10^7$, generated via a Sobol se-
812 quence in the eight-dimensional unit hypercube. The relative error on these
813 estimates (approximately equal to $[(p(1 - p)/N) / p$, where p is the estimated
814 fraction) is therefore order of 10^{-2} and 10^{-3} in the two cases respectively.
815 This makes both estimates very accurate.

816 In this case, we do not need to proceed to further waves. For all the
817 months, the emulator uncertainty (Table 4) is already 1-2 orders of magnitude
818 lower than the combined one from MD and OE. Additional waves therefore
819 leave the implausibility measure essentially unchanged. We comment further
820 on this in Section 6.

821 *5.4.1. Visualisation of the Non-Implausible Region*

822 The non-imausible region lives in an eight-dimensional space, one di-
823 mension per input parameter. To identify variables which play a role in
824 constraining the region, we look at two-dimensional scatter plots of the non-
825 implausible points, for all possible pairs of the 8 variables. For brevity, we do
826 not show these all here. However, in the 10% MD case, the plots reveal that
827 two variables are particularly significant: infiltration rate (V_6) and cooking
828 energy consumption (V_8). One additional variable (heating setpoint, V_1) also
829 seem to play a role in identifying non-imausible points.

830 Fig. 5 contains minimum-implausibility (MI) and optical-depth (OD)
831 plots for all pairs of the three above variables. The MI plot of a pair of
832 parameters shows the minimum value of the implausibility measure $I(\mathbf{x})$
833 over the hidden six dimensions. Thus, values greater than $T = 4$ in a MI
834 plot identify pairs of the two parameters in question that will always be clas-
835 sified as implausible, irrespectively of the value taken by the remaining six
836 parameters. For values smaller than $T = 4$, the OD plots instead show the
837 fraction of non-imausible points, in the hypervolume extending over the six
838 hidden dimensions.

839 Fig. 5 reveals that only values of the gas cooking higher than 5% of to-
840 tal household energy demand are able to yield a match with the observed
841 monthly consumptions. Moreover, within the explored ranges, such values
842 should be paired with values of the infiltration rate (right panels) approx-
843 imately between 0.25 and 0.65 ach **uniformity in lower/capital ach. Mo-**
844 **hammad, preference?**). Note that the two variables seem to be relatively

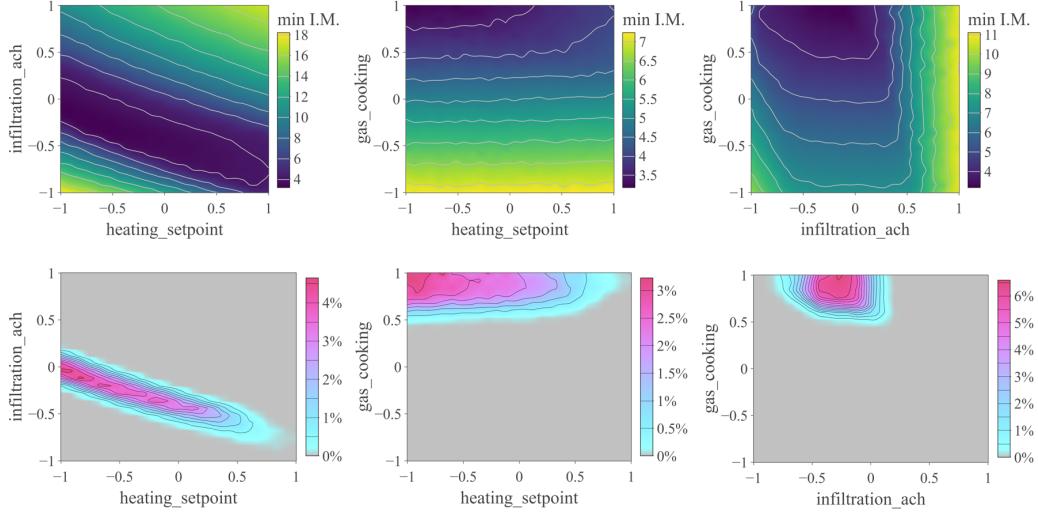


Figure 5: Plots of the non-imausible space when 10% MD is considered. Each panel in the top row shows the minimum implausibility along the six dimensions not included in the plot. Panels in the bottom row show the percentage of non-imausible region along the hidden dimensions.

845 independent. A stronger dependence can be seen instead for non-imausible
 846 values of heating setpoint and infiltration rate (left panels). Heating set-
 847 point values are by themselves only minimally constrained: however, higher
 848 infiltration rates yield more constrained, lower values of the heating setpoint.

849 The left subplots in Fig. 5 also suggest that non-imausible values of
 850 the heating setpoint may as well be found to the left of the range originally
 851 deemed appropriate for this input variable. A similar consideration is valid
 852 for gas cooking values higher than the relevant range. Similar findings are
 853 not uncommon, especially during a first wave of HM. In such a case, it is
 854 advisable to step back and expand the ranges in question, running addi-
 855 tional simulations in the new region before proceeding with HM. However,
 856 for brevity purposes and in accordance with the methodological aim of the
 857 work, in this illustrative example we limit our discussion to the hypercube
 858 with ranges specified in Table 1.

859 5.4.2. Sensitivity to Model Discrepancy Magnitude

860 The plots in Fig. 5 concern the case where the magnitude of model dis-
 861 crepancy is 10%. Fig. 6 shows MI and OD plots for the same three variables,
 862 in the case of 20% MD. Roughly, similar patterns to the ones of Fig. 5 emerge,

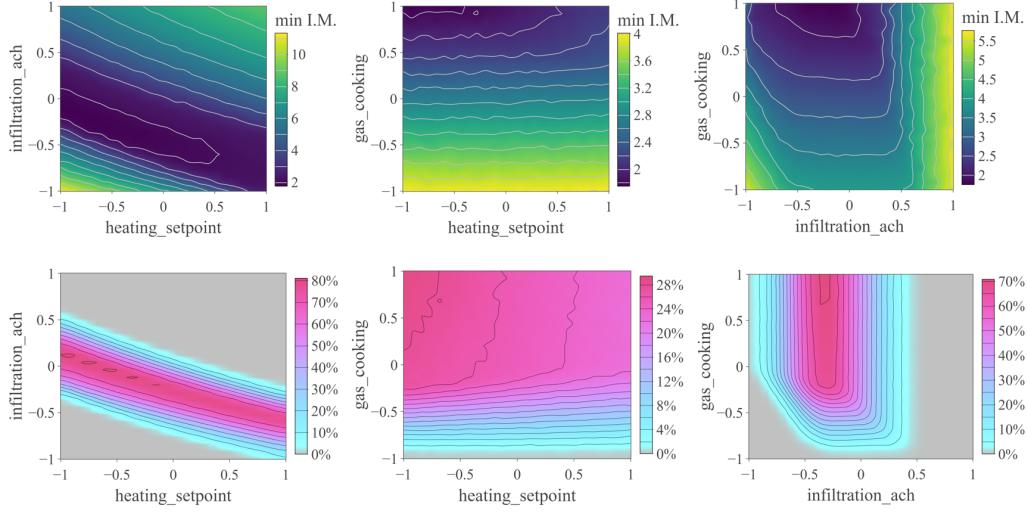


Figure 6: Same content as in Fig.4, when 20% MD is considered. Each panel in the top row shows the minimum implausibility along the six dimensions not included in the plot. Panels in the bottom row show the percentage of non-imausible region along the hidden dimensions.

albeit spread over larger regions: due to a lower confidence in the model, we rule out fewer points as implausible. The percentage of non-imausible points has risen to 19.52%, from only 0.3% in the 10% MD case.

This comparison highlights the potential sensitivity of history matching results to the choice of MD. Note that the precision of the emulator(s) used also plays a role in this. In our case, as Table 4 shows, we have remarkably precise emulators. Hence, essentially all the uncertainty accounted for in comparing emulator predictions and observations comes from model discrepancy and, to a lower extent, measurement error. Doubling the MD will thus make points significantly more likely to be deemed non-imausible. This once more highlights the importance of assessing the right order of magnitude of model discrepancy prior to performing history matching, as discussed in Sections 2.3 and 3.2.

5.4.3. Role of Different Constraints and Correlation among Outputs

The emulators allow to explore a wide range of questions, for which it would otherwise be impossible to draw sound inference. We briefly discuss one such example here, by looking at the compatibility between the imposed constraints (observed consumptions) and the correlation of model predictions

881 across different months. Note, in fact, that a strong correlation between two
882 model outputs restricts the pairs of observed consumptions that can lead to
883 a match.

884 Each of the upper-diagonal panels of Fig. 7 shows a scatter plot of emulated
885 gas consumption of the relevant pair of months, on a random sample
886 of 5,000 inputs. The observed consumption to be matched is identified by a
887 red cross. In most cases the latter is outside the region of outputs. However,
888 the presence of model discrepancy (10%) and observational error allows to
889 identify a region of non-implausible outputs, which is highlighted in turquoise
890 in each plot. Note, this is a subregion of the output space, not of the input
891 space as in Figures 5 and 6. A zoom of this region is displayed in the lower-
892 diagonal panels, with points coloured by the overall IM (9). The observations
893 aimed to be matched are displayed within a shaded rectangle identifying 5%
894 observational error.

895 Finally, the histograms on the diagonal of the same figure show the distribution
896 of the IM of a given month, on the points which would be classified as
897 non-implausible if that month was not considered. The IM is reported with
898 its original sign (expression (8) without absolute value), with positive values
899 denoting an emulated consumption higher than the observed one. Values
900 outside the range $[-4, 4]$ are highlighted by a darker colour. Months such as
901 January, March, November and December, where all values are between -4
902 and 4 , do not contribute to reduce the space once the other eight months'
903 constraints are included. On the other hand, a month such as September
904 rules out around 85% of the space that would be considered non-implausible
905 without it.

906 As the upper-diagonal panels reveal, some of the simulated monthly consumptions
907 are highly correlated, *e.g.*, January's and February's. In similar
908 cases, imposing that one month's observed consumption is matched will limit
909 which observed consumption of the other month can also be matched. In the
910 case of January and February, the location of the cross is very close to the line
911 of simulated consumptions, hence both observations can be easily matched
912 simultaneously (and can also be matched at the same time of all other seven
913 months, as we have seen). However, despite their high correlation, the two
914 months do not play an interchangeable role. The January constraint plays
915 a redundant role the February one is accounted for. But the opposite is
916 not true, as the x -range of the two diagonal histograms reveal. The same is
917 true for other pairs of strongly correlated outputs (*e.g.*, Feb-Mar, Feb-Nov,
918 Oct-Nov to a lower extent).

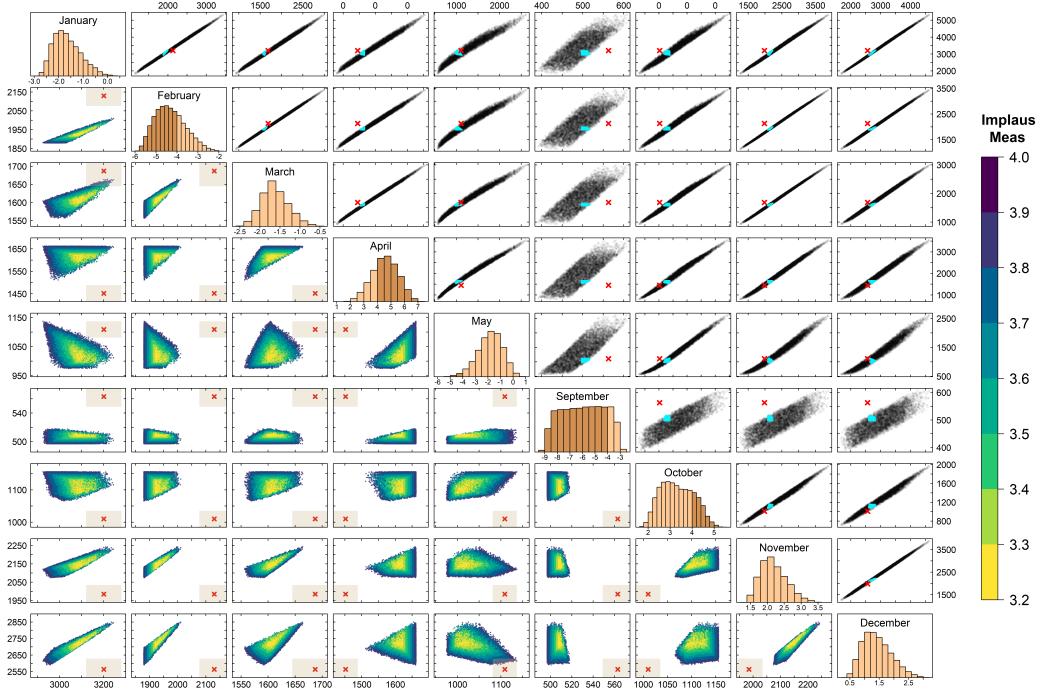


Figure 7: Information on model outputs, observations, and each month's contribution to cutting the non-implausible region. Upper-diagonal panels: scatter plot of emulated gas consumption for each pair of months (randomly selected 10,000 inputs). The red cross locates the observation to be matched, the light blue stain locates the non-implausible region. Lower-diagonal panels: zoom of the non-implausible region, coloured by implausibility measure (IM). Shaded rectangle around the cross identifies 5% measurement error. Diagonal panels: distribution of each month's IM, on the space deemed non-implausible when only constraints from the remaining months are considered. The darker colour denotes values outside the interval $[-4, 4]$, *i.e.* identifies inputs which transition from being non-implausible to being implausible when that month's constrain is added.

919 Finally, we notice that many of the lower diagonal projections appear
920 delimited by straight lines. These generally represent the “boundary” beyond
921 which one of the month’s IM exceeds 4 in absolute value. This is particularly
922 evident for the months whose diagonal histograms extend beyond the interval
923 $[-4, 4]$. In this case, the non-implausible region is “cut” perpendicularly to
924 that month’s axis, on the lower or higher side of the range according to
925 whether the histogram exceeds -4 or 4 .

926 *5.5. Adding Temperature Constraints*

927 Implausibility measures were previously used to history-match observed
928 gas consumption. They can however be used with any model output for which
929 observations are available, a feature that marks a difference with measures
930 such as MBE and CV(RMSE), as discussed in Section 2.3. In our case, we can
931 history match the model temperature predictions in the kitchen and master
932 bedroom, for which hourly time series are available both as observations (field
933 data) and as model outputs.

934 Accounting for uncertainty when comparing simulations and observations
935 is more challenging for time series than for scalar quantities. Appropriate
936 tools may involve accounting for correlation across time, dimension-reduction
937 techniques (*e.g.*, PCA) and multidimensional implausibility measures. It is
938 beyond the scope of this work to go into these details, some of which are active
939 areas of statistical research. However, scalar quantities may be extracted
940 from a time series and history-matched through the same methodology dis-
941 cussed in Section 3. We show an example here.

942 In order to include summer constraints, we consider the following scalar
943 quantity: the average temperature difference in July between day (8am–
944 11pm) and night (00am–7am) hours. For each of the two rooms, we compute
945 the above quantity across the 1,000 design simulations, build an emulator of
946 it as a function of the 8 input parameters, and history-match the parameters.

947 Notice that the choice of matching temperature difference introduces a
948 different type of constraint to the one imposed on gas, where absolute con-
949 sumption was matched. We consider the same levels of MD and OE as in
950 Section 5.3. Accounting for each of the two temperature constraints sep-
951 arately rules out as implausible 17.96% (kitchen) and 24.30% (master) of the
952 space respectively (we use here the threshold $T = 3$ since only one condition
953 was considered in each of the two cases). Unsurprisingly, we note a strong
954 dependence between the constraints coming from the two rooms: 12.38% of
955 the space is classified non-implausible with respect to both constraints at

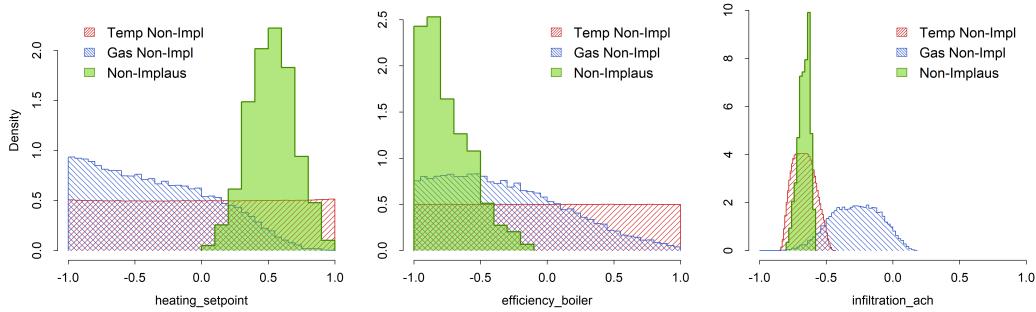


Figure 8: Marginal distribution of three input variables on non-implausible points. The shaded histograms mirror non-implausibility with respect to gas or temperature constraints only, the filled histogram with respect to both.

956 the same time, which is almost three times more than expected if they were
 957 independent.

958 Accounting for both energy (gas) and environmental (temperature) con-
 959 straints at the same time classifies only 0.00975% of the original eight-
 960 dimensional cube as non-implausible. This is about four times less than
 961 it would be expected if the two constraints were independent. The interplay
 962 between energy and environmental constraints is highlighted in Fig. 8, in
 963 the case of three input variables. For each one of them, the plot shows the
 964 distribution of non-implausible values for that variable, when the constraints
 965 from either gas, or temperature, or both are considered.

966 As illustration, we comment here on the left plot of Fig. 8, which concerns
 967 non-implausible values of the heating setpoint (HS) parameter.

968 On the one hand, if we only account for the July temperature con-
 969 straint, the original uniform distribution of heating setpoint values remains
 970 unchanged (red shade): this is expected, as heating is switched off in summer
 971 and does not, therefore, affect temperature. On the other hand, imposing
 972 only the gas constraint shifts the uniform distribution of HS towards the left.
 973 However, when both gas and temperature constraints are considered, only
 974 higher HS values turn out to be non-implausible. The explanation of this is
 975 as follows: third parameters are present whose values are in fact constrained
 976 by temperature, and which are correlated with HS. Forcing these parameters
 977 to have temperature-compatible values rules out low values of HS.

978 Something analogous happens with the boiler efficiency parameter. Fi-
 979 nally, the infiltration rate (which plays an important role in limiting temper-

980 ature oscillation in the building) is indeed mainly driven by the temperature
981 constraint. The additional imposition of the gas constraint only slightly re-
982 duces its non-imausible range.

983 6. Conclusion and Extensions of the Work

984 This work has discussed robust ways of accounting for uncertainty in pre-
985 calibration and calibration of building energy models. The methodology has
986 been illustrated on an actual dwelling and its energy model, with the aim
987 of simultaneously matching observations of different nature: energy (gas and
988 electricity use) and environmental (temperature in two zones).

989 Our methodology operates in a framework that allows to account for
990 different sources of quantified uncertainty. We have discussed typical sources,
991 such as model discrepancy and measurement error, and shown that their
992 magnitude can play a decisive role in assessing the goodness of a match
993 between simulations and observations. We have therefore illustrated History
994 Matching (HM), a statistical procedure that sequentially rules out regions of
995 the input space where simulations cannot match observations.

996 To explore the model’s response throughout the input space, HM makes
997 use of statistical emulators. Once built and validated, an emulator only
998 requires milliseconds on a personal laptop to predict the model response on
999 a new input configuration, and its prediction uncertainty can be naturally
1000 embedded in the HM framework. The hmer R package allows to implement
1001 emulation and the whole HM procedure [34].

1002 While HM generally proceeds in multiple waves, in the example study dis-
1003 cussed in this work one wave was enough to identify the final non-imausible
1004 region. This is a peculiarity of this example, where remarkably precise emu-
1005 lators for all quantities of interest could be built over the whole input space.
1006 In general, one of the strengths of HM is that the above is not at all required.
1007 Some of the outputs may not even be emulated in early stages. The sequential
1008 procedure allows the researcher to add constraints as the procedure moves to
1009 later waves. A quantity that was difficult to emulate in a given wave often
1010 behaves more smoothly within the more constrained NROY region of a later
1011 wave, and it can therefore be emulated with great precision at that stage.

1012 Also notice that in our case, the speed of the E+ model allowed us to
1013 run 1,000 simulations initially. Such a large number of runs is, once again,
1014 not needed in general. In fact, HM only needs a few runs at the start. The
1015 procedure itself will guide the researcher to focus new simulations within

1016 a sequence of smaller regions, that eventually converge to the final non-
1017 implausible region. In additiona experiments that has not been detailed in
1018 this work , we have randomly selected 100 of the original 1,000 simulations
1019 and used these to train (80 runs) and validate (20 runs) the new emulators.
1020 However, the quadratic behaviour of the simulator was easily captured by the
1021 80 simulations, yielding results that, after only one wave of HM, were prac-
1022 tically identical to the ones shown in Section 5, where all 1,000 simulations
1023 were used.

1024 Whilst emulators are classically used as model surrogates in HM, it is
1025 worth observing that less precise surrogates of the model may also be used
1026 (*e.g.*, linear regression), especially during the first explorative stages of re-
1027 search. Using linear regression, for example, would only take a few minutes of
1028 computational power but still allows the researcher to get a feeling of where
1029 the non-implausible region lies, while operating in a sound uncertainty frame-
1030 work.

1031 In terms of computational power, the efficiency of emulation over the
1032 original simulator is apparent. All emulation and HM computations carried
1033 out in this work were performed on a personal laptop with 16GB RAM and
1034 1.9GHz processor. For gas, and without any parallelisation, the computations
1035 (emulation and HM at 10^7 inputs, for nine months) were carried out in less
1036 than 28 hours. While still several orders of magnitude faster than the model
1037 running times, this is relatively slow for emulation standards, due to the high
1038 number of training points used (700). With only 80 training points, the same
1039 computations on the same device are executed in 2h and 5min. If we were
1040 to use the original E+ simulator directly, 10^7 simulations would have instead
1041 required about 18 years (as our original 1,000 simulations were performed in
1042 about 16 hrs).

1043 Finally, extensions of the methodology illustrated here can be manifold.
1044 Suppose a simulator is exceptionally slow to run, and faster versions produce
1045 lower-resolution results. A Bayesian framework where accurate emulators of
1046 the expensive model can be built, by using results of the faster version as
1047 prior information, is illustrated in [40]. This approach is referred to as multi-
1048 level (or multi-fidelity) emulation; an application for transport of volcanic
1049 ashes is discussed in [41]. Applications in the energy sector may concern the
1050 case of multi-vector energy modelling across districts or even national level,
1051 which is becoming crucial as a diverse mix of renewable inputs are used to
1052 replace fossil-based alternatives. Another similar example is building fabric
1053 energy modelling at building, cluster or even urban level whereby increasing

1054 the physical boundary and resolution of the model imposes a very heavy
1055 computational duty.

1056 An additional extension of the methodology we have discussed concerns
1057 the HM of time series data. This requires the use of emulators for time-
1058 evolving systems, as discussed for example in [42]. Multivariate implausibility
1059 measures ([14, 43]) may then be used to assess the distance between predicted
1060 and observed time series. This however requires a thoughtful specification
1061 of the time-correlation not only for the time series in question, but also for
1062 MD and OE. A discussion of this is beyond the scope of this work, but
1063 the interested researcher may refer to [ref] for an example of application to
1064 ***.(Michael, Hailiang, suggestion here?)

1065 References

- 1066 [1] S. Imam, D. A. Coley, I. Walker, The building performance gap: Are
1067 modellers literate?, *Building Services Engineering Research and Tech-*
1068 *nology* 38 (3) (2017) 351–375.
- 1069 [2] M. Royapoor, A. Antony, T. Roskilly, A review of building climate and
1070 plant controls, and a survey of industry perspectives, *Energy and Build-*
1071 *ings* 158 (2018) 453–465.
- 1072 [3] ASHRAE, Guideline 14-2002: Measurement of Energy and Demand Sav-
1073 *ings*, ASHRAE, Atlanta (2002).
- 1074 [4] M. Royapoor, T. Roskilly, Building model calibration using energy and
1075 environmental data, *Energy and buildings* 94 (2015) 109–120.
- 1076 [5] D. Hou, I. Hassan, L. Wang, Review on building energy model calibra-
1077 *tion by Bayesian inference*, *Renewable and Sustainable Energy Reviews*
1078 143 (2021) 110930. doi:10.1016/j.rser.2021.110930.
- 1079 [6] M. Goldstein, N. Huntley, Bayes Linear Emulation, History Matching,
1080 and Forecasting for Complex Computer Simulators, in: R. Ghanem,
1081 D. Higdon, H. Owhadi (Eds.), *Handbook of Uncertainty Quantifica-*
1082 *tion*, Springer International Publishing, 2017, pp. 9–32. doi:10.1007/
- 1083 978-3-319-12385-1_14.
- 1084 [7] M. C. Kennedy, A. O'Hagan, Bayesian Calibration of Computer Models,
1085 *Journal of the Royal Statistical Society: Series B* 63 (3) (2001) 425–464.

- 1086 [8] A. O'Hagan, Bayesian analysis of computer code outputs: A tutorial,
1087 Reliability Engineering & System Safety 91 (10-11) (2006) 1290–1300.
- 1088 [9] P. S. Craig, M. Goldstein, J. C. Rougier, A. H. Seheult, Bayesian Fore-
1089 casting for Complex Systems Using Computer Simulators, Journal of
1090 the American Statistical Association 96 (454) (2001) 717–729.
- 1091 [10] T. L. Edwards, S. Nowicki, B. Marzeion, R. Hock, H. Goelzer,
1092 H. Seroussi, et al., Projected land ice contributions to twenty-first-
1093 century sea level rise, Nature 593 (7857) (2021) 74–82.
- 1094 [11] A. L. Wilson, C. J. Dent, M. Goldstein, Quantifying uncertainty in
1095 wholesale electricity price projections using Bayesian emulation of a gen-
1096 eration investment model, Sustainable Energy, Grids and Networks 13
1097 (2018) 42–55.
- 1098 [12] I. Andrianakis, I. R. Vernon, N. McCreesh, T. J. McKinley, J. E. Oakley,
1099 R. N. Nsubuga, M. Goldstein, R. G. White, Bayesian history matching
1100 of complex infectious disease models using emulation: a tutorial and a
1101 case study on HIV in Uganda, PLoS computational biology 11 (1) (2015)
1102 e1003968.
- 1103 [13] P. S. Craig, M. Goldstein, A. H. Seheult, J. A. Smith, Pressure match-
1104 ing for hydrocarbon reservoirs: a case study in the use of Bayes linear
1105 strategies for large computer experiments, in: Case Studies in Bayesian
1106 Statistics, Vol. 121 of Lecture notes in statistics, Springer, 1997, pp.
1107 37–93. doi:10.1007/978-1-4612-2290-3_2.
- 1108 [14] I. Vernon, M. Goldstein, R. G. Bower, Galaxy formation: a Bayesian
1109 uncertainty analysis, Bayesian analysis 5 (4) (2010) 619–669. doi:10.
1110 1214/10-BA524.
- 1111 [15] D. Domingo, I. Malmierca-Vallet, L. Sime, J. Voss, E. Capron, Using ice
1112 cores and Gaussian process emulation to recover changes in the Green-
1113 land ice sheet during the last interglacial, Journal of Geophysical Re-
1114 search: Earth Surface 125 (5) (2020). doi:10.1029/2019JF005237.
- 1115 [16] I. Vernon, J. Liu, M. Goldstein, J. Rowe, J. Topping, K. Lindsey,
1116 Bayesian uncertainty analysis for complex systems biology models: emu-
1117 lation, global parameter searches and evaluation of gene functions, BMC
1118 systems biology 12 (1) (2018) 1–29.

- 1119 [17] D. Daly, P. Cooper, Z. Ma, Understanding the risks and uncertainties
1120 introduced by common assumptions in energy simulations for Australian
1121 commercial buildings, Energy and Buildings 75 (2014) 382–393. doi:
1122 [10.1016/j.enbuild.2014.02.028](https://doi.org/10.1016/j.enbuild.2014.02.028).
- 1123 [18] A. Rasooli, L. Itard, C. I. Ferreira, A response factor-based method for
1124 the rapid in-situ determination of wall's thermal resistance in existing
1125 buildings, Energy and Buildings 119 (2016) 51–61. doi:[10.1016/j.enbuild.2016.03.009](https://doi.org/10.1016/j.enbuild.2016.03.009).
- 1127 [19] A.-H. Deconinck, S. Roels, Comparison of characterisation methods
1128 determining the thermal resistance of building components from on-
1129 site measurements, Energy and Buildings 130 (2016) 309–320. doi:
1130 [10.1016/j.enbuild.2016.08.061](https://doi.org/10.1016/j.enbuild.2016.08.061).
- 1131 [20] X. Meng, B. Yan, Y. Gao, J. Wang, W. Zhang, E. Long, Factors affecting
1132 the in situ measurement accuracy of the wall heat transfer coefficient
1133 using the heat flow meter method, Energy and Buildings 86 (2015) 754–
1134 765. doi:[10.1016/j.enbuild.2014.11.005](https://doi.org/10.1016/j.enbuild.2014.11.005).
- 1135 [21] G. Ficco, F. Iannetta, E. Ianniello, F. R. d. Alfano, M. Dell'Isola, U-
1136 value in situ measurement for energy diagnosis of existing buildings,
1137 Energy and Buildings 104 (2015) 108–121. doi:[10.1016/j.enbuild.2015.06.071](https://doi.org/10.1016/j.enbuild.2015.06.071).
- 1139 [22] K. Gaspar, M. Casals, M. Gangolells, A comparison of standardized
1140 calculation methods for in situ measurements of façades U-value, Energy
1141 and Buildings 130 (2016) 592–599. doi:[10.1016/j.enbuild.2016.08.072](https://doi.org/10.1016/j.enbuild.2016.08.072).
- 1143 [23] G. Desogus, S. Mura, R. Ricciu, Comparing different approaches to in
1144 situ measurement of building components thermal resistance, Energy
1145 and Buildings 43 (10) (2011) 2613–2620. doi:[10.1016/j.enbuild.2011.05.025](https://doi.org/10.1016/j.enbuild.2011.05.025).
- 1147 [24] C. Hoffmann, A. Geissler, The prebound-effect in detail: real indoor
1148 temperatures in basements and measured versus calculated U-values,
1149 Energy Procedia 122 (2017) 32–37. doi:[10.1016/j.egypro.2017.07.301](https://doi.org/10.1016/j.egypro.2017.07.301).

- 1151 [25] A. Marshall, R. Fitton, W. Swan, D. Farmer, D. Johnston, M. Benjaber,
1152 Y. Ji, Domestic building fabric performance: Closing the gap between
1153 the in situ measured and modelled performance, Energy and Buildings
1154 150 (2017) 307–317. doi:10.1016/j.enbuild.2017.06.028.
- 1155 [26] P. Baker, U-values and traditional buildings, Historic Scotland Conserv-
1156 ation Group: Glasgow, UK (2011).
- 1157 [27] H. Yassaghi, N. Mostafavi, S. Hoque, Evaluation of current and fu-
1158 ture hourly weather data intended for building designs: A Philadel-
1159 phia case study, Energy and Buildings 199 (2019) 491–511. doi:
1160 10.1016/j.enbuild.2019.07.016.
- 1161 [28] T. Hong, S. C. Taylor-Lange, S. D’Oca, D. Yan, S. P. Corgnati, Ad-
1162 vances in research and applications of energy-related occupant be-
1163 havior in buildings, Energy and buildings 116 (2016) 694–702. doi:
1164 10.1016/j.enbuild.2015.11.052.
- 1165 [29] W. Tian, Y. Heo, P. De Wilde, Z. Li, D. Yan, C. S. Park, X. Feng,
1166 G. Augenbroe, A review of uncertainty analysis in building energy as-
1167 sessment, Renewable and Sustainable Energy Reviews 93 (2018) 285–
1168 301. doi:10.1016/j.rser.2018.05.029.
- 1169 [30] S. Coveney, R. H. Clayton, Fitting two human atrial cell models to ex-
1170 perimental data using Bayesian history matching, Progress in biophysics
1171 and molecular biology 139 (2018) 43–58. doi:10.1016/j.pbiomolbio.
1172 2018.08.001.
- 1173 [31] M. Goldstein, D. Wooff, Bayes Linear Statistics: Theory and Methods,
1174 Wiley Series in Probability and Statistics, John Wiley & Sons, 2007.
- 1175 [32] C. E. Rasmussen, C. K. Williams, Gaussian Processes for Machine
1176 Learning, The MIT Press, Massachusetts Institute of Technology, 2006.
- 1177 [33] F. Pukelsheim, The three sigma rule, The American Statistician 48 (2)
1178 (1994) 88–91. doi:10.2307/2684253.
- 1179 [34] A. Iskauskas, hmer: History Matching and Emulation Package (2022).
1180 URL <https://CRAN.R-project.org/package=hmer>

- 1181 [35] J. Kragh, J. Rose, H. N. Knudsen, O. M. Jensen, Possible explanations
1182 for the gap between calculated and measured energy consumption of new
1183 houses, Energy Procedia 132 (2017) 69–74. doi:10.1016/j.egypro.
1184 2017.09.638.
- 1185 [36] J. L. Loeppky, J. Sacks, W. J. Welch, Choosing the Sample Size of a
1186 Computer Experiment: A Practical Guide, Technometrics 51 (4) (2009)
1187 366–376. doi:10.1198/TECH.2009.08040.
- 1188 [37] N. S. Lord, M. Crucifix, D. J. Lunt, M. C. Thorne, N. Bounceur,
1189 H. Dowsett, C. L. O'Brien, A. Ridgwell, Emulation of long-term changes
1190 in global climate: application to the late Pliocene and future, Climate
1191 of the Past 13 (11) (2017) 1539–1571. doi:10.5194/cp-13-1539-2017.
- 1192 [38] C. A. Pope, J. P. Gosling, S. Barber, J. S. Johnson, T. Yamaguchi,
1193 G. Feingold, P. G. Blackwell, Gaussian Process Modeling of Hetero-
1194 geneity and Discontinuities Using Voronoi Tessellations, Technometrics
1195 63 (1) (2021) 53–63. doi:10.1080/00401706.2019.1692696.
- 1196 [39] M. D. McKay, R. J. Beckman, W. J. Conover, A Comparison of Three
1197 Methods for Selecting Values of Input Variables in the Analysis of
1198 Output From a Computer Code, Technometrics 42 (1) (2000) 55–61.
1199 doi:10.1080/00401706.2000.10485979.
- 1200 [40] J. A. Cumming, M. Goldstein, Small Sample Bayesian Designs for
1201 Complex High-Dimensional Models Based on Information Gained Us-
1202 ing Fast Approximations, Technometrics 51 (4) (2009) 377–388. doi:
1203 10.1198/TECH.2009.08015.
- 1204 [41] N. J. Harvey, N. Huntley, H. F. Dacre, M. Goldstein, D. Thomson,
1205 H. Webster, Multi-level emulation of a volcanic ash transport and dis-
1206 perssion model to quantify sensitivity to uncertain parameters, Nat-
1207 ural hazards and earth system sciences 18 (1) (2018) 41–63. doi:
1208 10.5194/nhess-18-41-2018.
- 1209 [42] S. Conti, J. P. Gosling, J. E. Oakley, A. O'Hagan, Gaussian process
1210 emulation of dynamic computer codes, Biometrika 96 (3) (2009) 663–
1211 676. doi:10.1093/biomet/asp028.
- 1212 [43] J. A. Cumming, M. Goldstein, Bayes Linear Uncertainty Analysis for Oil
1213 Reservoirs Based on Multiscale Computer Experiments, in: A. O'Hagan,

₁₂₁₄ M. West (Eds.), The Oxford Handbook of Applied Bayesian Analysis,
₁₂₁₅ Oxford University Press Oxford, 2010, pp. 241–270.