

1 **Process-based climate model development harnessing**
2 **machine learning: I. a calibration tool for**
3 **parameterization improvement**

4 **Fleur Couvreur¹, Frédéric Hourdin², Daniel Williamson^{3,5}, Romain Roehrig¹,**
5 **Victoria Volodina⁵, Najda Villefranque^{1,4}, Catherine Rio¹, Olivier Audouin¹,**
6 **James Salter^{3,5}, Eric Bazile¹, Florent Brient¹, Florence Favot¹, Rachel**
7 **Honnert¹, Marie-Pierre Lefebvre^{1,2}, Jean-Baptiste Madeleine², Quentin**
8 **Rodier¹, Wenzhe Xu³**

9 ¹CNRM, University of Toulouse, Meteo-France, CNRS, Toulouse, France

10 ²LMD-IPSL, Sorbonne University, CNRS, 4 pl Jussieu, Paris, France

11 ³Exeter University, Exeter, United Kingdom

12 ⁴LAPLACE, University of Toulouse, CNRS, Toulouse, France

13 ⁵The Alan Turing Institute, 96 Euston Road, London, United Kingdom

14 **Key Points:**

- 15 • We apply Uncertainty Quantification to Single-Column Model/LES comparison
16 to calibrate free parameters
- 17 • We revisit model development strategy with an emphasis on processes for model
18 calibration
- 19 • The proposed tuning tool allows to formalize the complementary use of multicases
20 with various metrics

Corresponding author: Fleur Couvreur, fleur.couvreur@meteo.fr

Abstract

The development of parameterizations is a major task in the development of weather and climate models. Model improvement has been slow in the past decades, due to the difficulty of encompassing key physical processes into parameterizations, but also of calibrating or ‘tuning’ the many free parameters involved in their formulation. Machine learning techniques have been recently used for speeding up the development process. While some studies propose to replace parameterizations by data-driven neural networks, we rather advocate that keeping physical parameterizations is key for the reliability of climate projections. In this paper we propose to harness machine learning to improve physical parameterizations. In particular we use Gaussian process-based methods from uncertainty quantification to calibrate the model free parameters at a process level. To achieve this, we focus on the comparison of single-column simulations and reference large-eddy simulations over multiple boundary-layer cases. Our method returns all values of the free parameters consistent with the references and any structural uncertainties, allowing a reduced domain of acceptable values to be considered when tuning the 3D global model. This tool allows to disentangle deficiencies due to poor parameter calibration from intrinsic limits rooted in the parameterization formulations. This paper describes the tool and the philosophy of tuning in single-column mode. Part 2 shows how the results from our process-based tuning can help in the 3D global model tuning.

1 Introduction

Atmospheric global or regional circulation models used either for numerical weather prediction (NWP) or climate studies encompass a dynamical core and a physical component. The dynamical core computes the spatio-temporal evolution of atmospheric state variables by solving a discrete version of the fluid dynamic equations. The physical component quantifies the impact on the resolved variables of radiative, thermodynamical and chemical processes, as well as dynamical processes that occur at scales smaller than the computational grid. These processes are handled by a suite of sub-models, most often referred to as parameterizations, which provide source terms in the resolved-scale equations. Parameterizations (e.g., turbulence, convection, radiation, microphysics) are often based on a mixture of physical principles and heuristic description of the involved processes, of their interactions and of their impact on the larger resolved scales. Although it is difficult to trace back the origin of the term “parameterization” in climate model-

ing, it semantically points to the fact that the sub-models summarize the processes as functions of the model state vector \mathbf{x} (typically the value of zonal and meridional wind, temperature and water phases at each point of the 3D model grid) that depends on some free parameters. These free parameters arise from the simplification of the complex nature of the subgrid processes (e.g., assuming a bulk thermal plume instead of a population of plumes, stationarity). The atmospheric model can be summarized as

$$\frac{\partial \mathbf{x}}{\partial t} = \mathcal{D}(\mathbf{x}) + \sum_p \mathcal{P}_p(\mathbf{x}, \boldsymbol{\lambda}_p) \quad (1)$$

41 where \mathcal{D} stands for the discretized form of the fluid dynamic equations, \mathcal{P}_p for the source
 42 term provided by the parameterization of the process p and $\boldsymbol{\lambda}_p$ for the associated free
 43 parameters. This equation may however be too simplistic, as, in reality, a given param-
 44 eterization often depends on intermediate variables provided by other parameterizations
 45 (e.g., cloud fraction used in radiation, turbulence variance used in the cloud scheme) and
 46 computes additional prognostic variables (e.g., turbulence kinetic energy). Nevertheless,
 47 with this simplified framework, improving models through parameterization development
 48 means both to propose more appropriate functional forms \mathcal{P}_p and to identify acceptable
 49 or better values of the free parameters $\boldsymbol{\lambda}_p$.

50 Among the different parameterizations, those involved in the representation of tur-
 51 bulence, convection and clouds still challenge state-of-the art NWP and climate mod-
 52 els (Holtslag et al., 2013; Nam et al., 2012; Nuijens et al., 2015; Klein et al., 2017; Ran-
 53 dall et al., 2003; Bony et al., 2015). Innovative and diverse concepts and ideas have been
 54 proposed over the past decade to improve this representation (Rio et al., 2019). A de-
 55 tailed understanding of the physical processes leading to the formation of low-level clouds
 56 can be obtained by Large-Eddy Simulations (LES) (Guichard & Couvreur, 2017), which
 57 reproduce, with high fidelity, the turbulent dynamics within the clouds (e.g., Siebesma
 58 & Cuijpers, 1995; Neggers, Duynkerke, & Rodts, 2003; Wang & Feingold, 2009). LES
 59 are therefore increasingly used to derive and evaluate the conceptual models at the root
 60 of boundary-layer and shallow cloud parameterizations. The choice of the parameter-
 61 ization free parameters is also crucial for the simulation of clouds. Their calibration or
 62 “tuning” consists in searching for acceptable or optimal values of these parameters, such
 63 that the associated model configuration has a realistic behavior under various conditions
 64 and compared to a suite of observations (Mauritsen et al., 2012). Calibration is there-
 65 fore a fundamental aspect of NWP or climate model development (Bellprat et al., 2012;

66 Duan et al., 2017; Schmidt et al., 2017). However, it is often conducted without much
67 control on the way it modifies the parameterization behavior at the process level as the
68 calibration focuses more on regional or global constraints, such as the radiative balance
69 of the Earth System for climate models, or performance metrics (e.g. root mean square
70 error, skill scores) for NWP models. Hourdin et al. (2017) compile the tuning strategies
71 of several climate groups and emphasize that most of the parameters used to tune cli-
72 mate models (droplet size, fall velocity, entrainment rate) are related to clouds (see also
73 Golaz et al., 2013), i.e. the most uncertain processes that affect radiation, the primary
74 engine of the atmospheric circulation.

75 Given the societal needs for reliable climate simulations and weather forecasts, the
76 progress achieved by the global atmosphere modeling community has been found slow
77 (Jakob, 2010). Several systematic errors in state-of-the-art models have been modestly
78 reduced, such as those regarding the surface temperature over the eastern oceans (Richter,
79 2015), the rainfall distribution in the Tropics (Flato et al., 2013), the variability of the
80 liquid water path (Jiang et al., 2012) and the low clouds (Nam et al., 2012). The dead-
81 lock of the cloud parameterization, highlighted by Randall et al. (2003), is still an issue
82 today. This too slow improvement of models can be attributed to remaining deficiencies
83 in the structure of the parameterization itself (the function \mathcal{P}_p) but also to the calibra-
84 tion of model parameters that can be considered as a bottleneck in model development.
85 On the one hand, the calibration may not be done efficiently enough, and on the other
86 hand, tuning may induce error compensations that contribute to slow model develop-
87 ment. Indeed, a new model development usually starts with a model score degradation
88 by breaking this compensation, as often experienced in the weather prediction centers
89 where strong weight on well-established metrics slows down the implementation of new
90 model development in the operational version (Sandu et al., 2013).

91 Various avenues have been proposed to get around these difficulties and acceler-
92 ate climate model improvement. A first avenue seeks to exploit the high resolution, ex-
93 plicitly resolving convection, to reduce the number of involved parameterizations. With
94 the recent increase of computer power, it is nowadays possible to run global kilometer-
95 scale resolution simulations over a few months (Satoh et al., 2008, 2019; Stevens et al.,
96 2019). However, the explicit simulation of the fluid dynamics associated with the life cy-
97 cle of a cumulus requires grid resolution of the order of several tens of meters. Such res-
98 olution will not be accessible in the foreseeable future for climate change projections which

99 require simulations of the global Earth System covering at least several hundreds of years
100 (model spin-up plus transient simulations in response to anthropogenic forcing). The super-
101 parameterization approach (Randall et al., 2003) proposes an intermediate pathway by
102 introducing a convection-permitting model in each column of a conventional general cir-
103 culation model (GCM) to replace the deep convection parameterization (Khairoutdinov
104 et al., 2005). The use of a large-eddy model instead of a convection-permitting model
105 in such framework further removes the boundary-layer and shallow convection param-
106 eterizations (Grabowski, 2016; Parishani et al., 2017). A second avenue recently explored
107 the potential of machine learning approaches, which ultimately envisions to replace some
108 parameterizations by neural networks or similar algorithms, properly trained on convection-
109 permitting model simulations or superparameterized GCM (Krasnopolsky et al., 2013;
110 Brenowitz & Bretherton, 2018; Gentine et al., 2018).

111 A third proposition consists in retaining parameterizations in models but adjoining
112 new tools relying on machine learning to accelerate model development. This choice
113 is motivated by the fact that parameterizations summarize our current understanding
114 of the dynamics and physics of atmospheric processes and offer the power of interpre-
115 tation, crucial to build our confidence in the extrapolation beyond observed conditions
116 realized by any climate projections. The ESM2.0, proposed by Schneider et al. (2017),
117 belongs to this category. The authors defend that the major progress in Earth-System
118 model development should come from a more systematic use of global observations and
119 high-resolution simulations thanks to machine learning algorithms. They also underline
120 the importance of climate model calibration. In particular, they stress that their new
121 Earth System modeling framework comes with challenges such as developing innovative
122 learning algorithms, identifying the best metrics, combining information from observa-
123 tions and high-resolution, innovating in the design of parameterizations to more easily
124 benefit from new observations or evolution of the models (e.g., refinement of resolution).

125 Along the same lines, we propose, in this paper, a new approach which allows the
126 development of the parametrizations and their calibration to be tackled at the same time.
127 We argue that a major slowdown of model improvement resides in the difficulty to clearly
128 identify parameterization deficiencies and to properly disentangle them from the inher-
129 ent calibration of their adjustable parameters at the process and global scales. It is likely
130 that process-scale parameterization improvements are often hidden by the unavoidable
131 full model re-tuning, required to maintain a reasonable radiative balance or acceptable

132 scores. In the proposed approach, machine learning is harnessed in a principled way to
 133 calibrate parameterizations at process level. We promote a more systematic use of the
 134 multi-case comparison between Single-Column Model (SCM) and LES to evaluate and
 135 calibrate parameterizations. Such a systematic use is not feasible however without more
 136 objective and automatic methods than the traditional trial/error approach used to fix
 137 parameter values during the parameterization development. Indeed, this trial/error ap-
 138 proach is only applicable to one piece of a particular parameterization and one or two
 139 relevant cases at most. Here, we aim at assessing a set of parameterizations \mathcal{P}_p for a se-
 140 ries of test cases, which can be formalized as the question of the existence of a sub-space
 141 of the parameters λ_p that allows to match metrics between SCM and LES results for the
 142 series of cases, within a given tolerance to error.

143 Hourdin et al. (2017) reviewed the general practice for climate model calibration
 144 and proposed three different levels of calibration in a model development: a first cali-
 145 bration at the level of individual parameterizations, then a calibration of each compo-
 146 nent of the Earth System model and eventually a calibration of the full Earth System
 147 model. Distinguishing those three levels may avoid compensating errors that could arise
 148 if the calibration is only done at the last level. In this paper, we propose a methodol-
 149 ogy to address the first phase, *i.e.* the process-level calibration and defend that it can
 150 be part of the elaboration of a well-defined calibration strategy based on solid physical
 151 and statistical methodologies. By doing so, we tackle model development and param-
 152 eter calibration together rather than independently as currently done for most climate
 153 model development.

154 Machine learning has already been proposed to calibrate free parameters (e.g., en-
 155 semble Kalman filters as in Schneider et al., 2017). The methodology retained here for
 156 model calibration uses history matching with Gaussian processes. History matching is
 157 an efficient way to explore and reduce the domain of free parameters λ_p and document
 158 how a model physics, namely the suite of functions \mathcal{P}_p , behaves within this domain. Williamson
 159 et al. (2013) applied history matching to tune the Hadley Climate Model and stressed
 160 its advantage: it accounts for the various sources of uncertainties in assessing the com-
 161 patibility of the model with the reference: namely the reference uncertainty itself, the
 162 uncertainty introduced by the Gaussian process representation of the parameterization,
 163 and the intrinsic ability of the model to represent the reference (often referred to as struc-
 164 tural error or model discrepancy). History matching inherently deals with the overcon-

165 fidence issue, which emerges when model calibration is addressed as an optimization prob-
166 lem (Salter et al., 2019). It has been widely used to calibrate models in astrophysics (Vernon
167 et al., 2010), epidemiology (Andrianakis et al., 2017) and hydrocarbon reservoirs (Craig
168 et al., 1996). It has been applied to climate models (Williamson et al., 2015, 2017) and
169 is starting to be used to find biases in models (McNeall et al., 2019).

170 Whilst history matching has been applied to calibrate 3D models, it has not been
171 harnessed for process-level tuning, as we advocate here through application to SCM/LES
172 comparison. The SCM approach provides confidence in the model’s ability to represent
173 some of the key processes whereas a direct calibration of the 3D global model targeting
174 large-scale constraints may hide compensating errors (as discussed in Williamson et al.,
175 2017). SCM calibration is able to reduce the domain of the free parameters for a param-
176 eterization, information that can be used for efficiently calibrating the full 3D global model
177 (as we demonstrate in part II). The breakthrough proposed here was only possible thanks
178 to a strong collaboration between the Uncertainty Quantification community and the at-
179 mospheric modelers.

180 The present paper focuses on parameterizations involved in the representation of
181 boundary-layer clouds. Indeed, well-established case studies exist for such regimes and
182 LES have been shown to realistically represent the main processes. However, this method-
183 ology can be easily expanded to other parameterizations and other objectives in the Earth
184 System.

185 The paper is organized as follows: the next section describes the SCM/LES frame-
186 work highlighting its advantages, recalls the different steps used in the development of
187 a parameterization and details the new philosophy advocated here. Section 3 presents
188 the statistical tool, with a focus on its philosophy and its main ingredients. Section 4
189 presents a guideline for its use based on a simple illustration. The paper ends with con-
190 clusions in section 5. A companion paper (part II) illustrates the significant advances
191 in model development offered by this tool. It exploits process-based calibration for model
192 development and shows how this tool provides guidance for the tuning of a 3D global
193 model.

194 **2 A systematic use of the SCM/LES comparison**

195 Although observations, especially combinations of observations, nowadays provide
 196 detailed information at high temporal and spatial resolution on the characteristics of con-
 197 vection and clouds (Masunaga, 2012; Kumar et al., 2015; Bouniol et al., 2016; Masunaga
 198 & Luo, 2016), their use for process-level analysis is still hampered by the difficulty of (i)
 199 comparing model output to what the satellite measurements exactly sample (although
 200 the model to satellite approach with simulators partly resolves this issue) and (ii) iden-
 201 tifying the physical processes responsible for such characteristics. Here, we promote the
 202 use of Large-Eddy Simulations for the following reasons. LES have the advantage of pro-
 203 viding coherent 3D fields characterizing the dynamical and thermodynamical state of the
 204 atmosphere. Of course, LES models include turbulence and microphysics parameteri-
 205 zations and thus contain modeling uncertainties, but they have been shown to reproduce
 206 the turbulent dynamics of the clouds with high fidelity (e.g., Neggers, Duynkerke, & Rodts,
 207 2003; Heus et al., 2009). As a result, LES have become a central tool in the development
 208 of parameterizations of convection and clouds. Their analysis has helped in building the
 209 conceptual models behind several parameterizations (e.g., Neggers et al., 2002; Rio et
 210 al., 2010). LES are also used for the evaluation of the parameterizations in particular
 211 those involved in the representation of boundary layers and shallow clouds (e.g., Ay-
 212 otte et al., 1996; Golaz et al., 2002; Hourdin et al., 2002; Neggers et al., 2004; Siebesma
 213 et al., 2007; Rio & Hourdin, 2008; Caldwell & Bretherton, 2009; Neggers, 2009; Pergaud
 214 et al., 2009; Rio et al., 2010; Suselj et al., 2013; Neggers et al., 2017; Tan et al., 2018;
 215 Suselj et al., 2019).

For their evaluation, parameterizations are often tested in a single-column frame-
 work, particularly relevant for global circulation model parameterizations, which are fun-
 damentally 1D. SCMs are built by extracting, from a 3D model, a single atmospheric
 column, which integrates the same set of subgrid parameterizations (boundary-layer, shal-
 low convection, deep convection and microphysics schemes) and is run in a constrained
 large-scale environment (Zhang et al., 2016). The state vector of the SCM simulation
 is then a restriction to one column \mathbf{x}_c of the full 3D state vector \mathbf{x} and Eq. 1 reduces
 to Eq. 2. The dynamical term $\mathcal{D}(\mathbf{x})$ becomes a source term \mathcal{F}_c specified as a function
 of time and altitude z ; we however discard this dependency in the notation for simplic-
 ity. It can also depend on the column full state vector, $\mathcal{F}_c(\mathbf{x}_c)$, if for instance the large-
 scale advection is separated between a prescribed horizontal advection and a vertical ad-

vection computed as $-w\partial\mathbf{x}_c/\partial z$, where w is an imposed vertical velocity. During the SCM integration, some parameterizations can be deactivated in which case the corresponding source term is either neglected or included in the forcing \mathcal{F}_c . It is the case for instance when the radiative heating is imposed rather than being computed interactively by the model radiation scheme or when turbulent surface fluxes are imposed rather than computed by the model bulk parameterizations. What really matters in the SCM/LES approach is that both models use the exact same initial and boundary conditions and forcing terms. In a simplified formalism, the SCM thus corresponds to

$$\frac{\partial\mathbf{x}_c}{\partial t} = \sum_{p \in \mathcal{P}_{\text{activated}}} \mathcal{P}_p(\mathbf{x}_c, \lambda_p) + \mathcal{F}_c(\mathbf{x}_c) \quad (2)$$

and the LES to

$$\frac{\partial\mathbf{y}}{\partial t} = \mathcal{L}(\mathbf{y}) + \mathcal{F}_c^*(\bar{\mathbf{y}}) \quad (3)$$

with

$$\mathbf{x}_c(t=0) = \bar{\mathbf{y}}(t=0) \quad (4)$$

216 where \mathbf{y} stands for the full LES state vector, $\mathcal{L}(\mathbf{y})$ to the LES model equations (which
 217 include the LES parameterizations), $\bar{\mathbf{y}}$ to the horizontal-domain average of the LES state
 218 vector and \mathcal{F}_c^* provides a 3D field but consists of the same forcing as the SCM, \mathcal{F}_c ap-
 219 plied identically on each individual column of the LES. The SCM/LES framework thus
 220 provides a rigorous comparison between both simulations, as it removes the uncertain-
 221 ties, which may arise from different initial conditions or large-scale forcing when directly
 222 comparing SCM to observations. This constrained framework also avoids the need to dis-
 223 entangle parameterization contributions from their coupling with the large-scale dynam-
 224 ics. Another important aspect of the method is that SCM simulations are computation-
 225 ally very cheap. The joint utilization of LES and SCM was first advocated by Randall
 226 et al. (1996); Ayotte et al. (1996) and has been, since then, widely used within the Global
 227 Energy and Water Exchanges (GEWEX) Cloud System Study (GCSS; Browning et al.
 228 (1993) community, now renamed the Global Atmospheric System Studies, GASS, com-
 229 munity). One of the most important legacies of this group for the atmospheric model-
 230 ing community is an ensemble of test cases that connect observations, LES and SCM,
 231 and which sample many typical situations over the globe, thought to be of importance
 232 for the climate system (e.g., Siebesma & Cuijpers, 1995; Brown et al., 2002; Duynkerke
 233 et al., 2004). As such, this framework has been increasingly used in model development
 234 (e.g., Hourdin et al., 2013; Gettelman et al., 2019; Hourdin et al., 2020; Roehrig et al.,

235 2020), all the more so as SCM simulations have been shown to reproduce uniquely the
236 behavior of their GCM justifying the use of SCM simulations for improving weather and
237 climate models (Hourdin et al., 2013; Neggers, 2015; Gettelman et al., 2019).

238 Traditionally, parameterizations are often tested over a few specific cases for which
239 high-resolution simulations are available (e.g., Ayotte et al., 1996). Recently, the im-
240 portance of using a wide benchmark of cases covering the different regimes encountered
241 in reality instead of only a limited number of cases has been stressed (e.g., Neggers et
242 al., 2012). We also highlight here the importance of using an extensive ensemble of cases.
243 The use of multi-case is indeed essential for exploring the various degrees of freedom of
244 the parameterization package. A stable boundary-layer case will constrain the turbulent
245 diffusion; the combination of cloud free and cumulus topped convective boundary lay-
246 ers will ensure that cloud cover is obtained for a good representation of convection; tran-
247 sition cases from stratocumulus to cumulus will ensure the extension to stratocumulus
248 regimes, etc. Combining multi cases and multi metrics is a much more robust assessment
249 of model performance as also highlighted by Neggers et al. (2017). To better use multi-
250 cases, one important technical aspect is a common definition, in a predefined acknowl-
251 edged format, for the description of the setup of reference cases, to be used both to per-
252 form SCM simulations or LES. This definition should include the description of the ini-
253 tial profiles and large-scale forcing but also contain information on the configuration to
254 be used (e.g. the type of surface boundary conditions, the existence of any nudging to-
255 wards reference vertical profiles, the way large-scale forcing are provided). An interna-
256 tional initiative is ongoing to agree on the description of the format for this definition
257 file. Such a standard format to define cases will ease the realization of cases by any model
258 and facilitate the share of new cases. The importance of creating libraries of high-resolution
259 simulations representing different climate is another important aspect already identified
260 as a goal by the GCSS community and stressed in Schneider et al. (2017). A common
261 format and the libraries of LES are an important pre-requisite for the tool presented here.
262 In addition, both will contribute to bringing the process-scale community and the com-
263 munity developing global models more closely together.

264 When comparing SCM and LES, the modeler has to decide which metrics to con-
265 sider. Various types of metrics can be used. One can directly compare components of
266 the SCM state vector \mathbf{x}_c to their equivalent in LES, the horizontal domain-average state
267 vector $\bar{\mathbf{y}}$ (e.g., vertical profiles of potential temperature, specific humidity and less of-

268 ten wind components). Assessing the ability of the parameterizations to reproduce the
269 time evolution of \mathbf{x}_c for a given forcing is indeed the ultimate goal. By doing so, one not
270 only tests the behavior of one particular parameterization but also its coupling with the
271 other parameterizations activated in the SCM. This may make the determination of the
272 behavior of the targeted parameterization more difficult and can hide compensating er-
273 rors: for example, a given temperature turbulent flux can be obtained by different con-
274 tributions from organized structures and small-scale turbulence when represented by two
275 different parameterizations such as in the Eddy-Diffusivity Mass-Flux framework (Hourdin
276 et al., 2002; Siebesma et al., 2007; Neggers, 2009; Pergaud et al., 2009). Another type
277 of metrics targets parameterization-oriented variables, such as mass fluxes, heating source
278 associated with one part of the motion only, subgrid-scale distribution of temperature
279 or water, cloud vertical structure, updraft vertical velocity, area fraction or entrainment
280 and detrainment rates. The metric, from the SCM point-of-view, is no-longer derived
281 from the model state variables but corresponds to a variable internal to the parameter-
282 izations. However, additional uncertainty arises from the way such variables and asso-
283 ciated metrics can be derived from LES. For example, clouds can be characterized in an
284 LES as all the grid cells containing condensed water (e.g., Siebesma & Cuijpers, 1995).
285 Combined with thresholds on the vertical velocity, cloudy updrafts can be separated from
286 cloudy downdrafts. The analysis of the joint distribution of variables or the use of ad-
287 hoc passive tracers can also be used in the LES to identify objects relevant with the con-
288 ceptual model of the parameterization (e.g., Couvreux et al., 2010; Rio et al., 2010; Chinita
289 et al., 2018; Brient et al., 2019). Such parameterization-oriented diagnostics have helped
290 in the refinement of the conceptual model at the root of the parameterization (e.g., Rio
291 et al., 2010; Jam et al., 2013; Rochetin et al., 2014). However, a question arises if such
292 diagnostics should also be used as metrics in the calibration process. Answering this ques-
293 tion on the relative importance to give to one type of metrics or another requires effi-
294 cient algorithms, as the one proposed here, to explore the various options. Note also that
295 using state vector-based metrics on a large set of cases that are more or less sensitive to
296 one aspect of the parameterization may help avoid the error compensation issue.

297 In line with Neggers et al. (2012), we advocate that, although not a new approach,
298 the power of SCM/LES comparisons is largely underestimated and under-exploited. Ap-
299 plying history matching to this comparison is a way to fully take advantage of the SCM/LES
300 on a large multi-case ensemble and explore whether there exists a sub-space of the pa-

parameter space for which the SCM is able to reproduce a series of LES simulations within a given uncertainty. Note that the metrics can be different from one case to the other. This tool offers the possibility to revisit the different intercomparison exercises documented in the literature and to benefit from this rich database still underused.

Eventually, a point that becomes crucial when using LES for parameterization evaluation and tuning is the assessment of LES reliability and its uncertainties. Although it has been shown, through the comparison to observations, that LES is able to correctly reproduce boundary-layer processes and shallow clouds (Couvreur et al., 2005; Neggers, Jonker, & Siebesma, 2003; Heus & Jonker, 2008), LES, as in many models, come with uncertainties associated to the advection scheme and the parameterizations still active in such simulations concerning small-scale turbulence, microphysics, radiation and surface fluxes. Sullivan and Patton (2011) have shown that a horizontal resolution of a few tens of meters for convective boundary layers is enough to get convergence for the mean, fluxes and variances but 10m resolution is needed in order to get convergence on skewness. The sensitivity of LES of shallow convection to resolution, size of the domain, sub-grid model and advection scheme has been widely investigated (Brown, 1999; Matheou et al., 2011; Pressel et al., 2017; Zhang et al., 2017; Wurps et al., 2020). In particular, it has been shown that most of the ensemble-averaged turbulence statistics are reasonably insensitive, allowing one to use LES results to develop and evaluate convection parameterizations. However, some characteristics of the cloud fields (e.g. size distribution of individual clouds) are more sensitive to resolution, advection scheme or subgrid-scheme (Brown, 1999; vanZanten et al., 2011; Pressel et al., 2017). For example, LES at 5-10m vertical resolution still have large uncertainties in boundary-layer regimes with sharp inversions where the LES subgrid turbulence parameterization is significantly active. Uncertainty around this reference should be documented so that history matching can explicitly take it into account.

3 *High-Tune Explorer* (htexplo), a statistical tool to calibrate model parameters and more

3.1 Overview

The present section describes the tool proposed to perform process-based calibration. Its objective is twofold: (i) characterize the domain of the model parameter values that allows the model to appropriately capture process-level metrics and which can

333 be used for subsequent calibration of the global model, and (ii) identify the model pa-
 334 rameters that limit model performance and thus highlight the need for model param-
 335 eterization revision. The tool relies on history matching approach developed by Vernon
 336 et al. (2010) and first used for climate studies by Williamson et al. (2013). This method
 337 aims at removing “unphysical” regions of parameter space iteratively, refocusing the search
 338 for “acceptably tuned” models at each step. The tool finds the subspace of the model
 339 parameter space containing simulations consistent with the reference metrics, acknowl-
 340 edging the various sources of uncertainty. This tool has already been successfully applied
 341 to identify the acceptable range of model parameter values in the 3D configuration of
 342 the Hadley Centre climate model (Williamson et al., 2013, 2015) or in the NEMO oceanic
 343 model (Williamson et al., 2017). It is here used for the first time in the context of the
 344 SCM/LES comparison for a given set of cases.

345 As already stated in the previous section, we focus here on the parameterizations
 346 involved in the representation of boundary-layer clouds (turbulence, convection, cloud
 347 micro and macrophysics, radiation). However, this methodology can be easily expanded
 348 to other parameterizations and other objects of the Earth system as soon as reliable ref-
 349 erences are available.

350 Figure 1 sketches the main steps of the *High-Tune Explorer* (htexplo in the follow-
 351 ing for an explorer to use High-resolution simulation to improve and Tune parameter-
 352 izations) tool:

- 353 • **1. Metric selection and references** First, the cases and associated target met-
 354 rics are selected. The relevant reference for each metric is then identified and the
 355 associated uncertainty is estimated. In the present case, the reference is an LES
 356 and the associated uncertainty is based on an LES ensemble. Observations could
 357 also be used with an associated error when an LES is not available. This phase
 358 is not model-specific and could be shared between different models.
- 359 • **2. Selection of model parameters** The model parameters to be calibrated are
 360 identified and their possible range of values are determined.
- 361 • **3. Experimental design and SCM runs** The experimental design consists of
 362 defining the ensemble of experiments (or SCM) to be run. The goal is to optimally
 363 sample the parameter space and provide a small set of parameter values for which

- 364 the single-column model will be run. Metrics are computed from each of the SCM
 365 simulations and form the training data-set on which emulators are built.
- 366 • **4. Building emulators**, i.e. construction of surrogate models, also called “em-
 367 ulators”, one for each metric. Each emulator is based on a Gaussian Process (GP)
 368 and predicts the corresponding metric value at any point of the full parameter space,
 369 without running the SCM. The GP statistical model also provides a probability
 370 distribution of its prediction, thus quantifying the prediction uncertainty for use
 371 in calibration.
 - 372 • **5. History matching** The comparison between the reference metrics and those
 373 inferred with the emulators is based on a distance that accounts for reference un-
 374 certainty, modeler tolerance to error or model discrepancy (induced by e.g., mis-
 375 representation of specific processes, inaccuracy of numerical solvers, model reso-
 376 lution) and emulator uncertainty. History matching rejects parameter values that
 377 lead to unacceptable model behavior (too large distance from the reference) and
 378 thus defines a not-ruled out yet (NROY) space, the model parameter space that
 379 cannot be further reduced given the sources of uncertainty.
 - 380 • **6. Iterative refocusing** To reduce the emulator uncertainty, but only where needed,
 381 new iterations (or waves) following steps 3 to 5 are performed, sampling the NROY
 382 space obtained at the end of the previous wave for the design and only construct-
 383 ing emulators over the NROY domain.

384 This tool is available freely under: <https://svn.lmd.jussieu.fr/HighTune>. Details on the
 385 different steps are given below. For simplicity, we first describe them for the first iter-
 386 ation and only one metric. Subsequent iterations and the addition of other metrics are
 387 discussed in section 3.7. This section ends with a discussion about the relationship be-
 388 tween the present tool and more common tools used for calibration and sensitivity anal-
 389 ysis.

390 **3.2 Step 1: Metric selection and references**

391 The metrics used to evaluate the SCM behavior depend on the physical situation
 392 considered and the parameterization hypothesis. Scalar metrics based on a dynamical
 393 or thermodynamical variable (e.g., potential temperature, water vapor mixing ratio, wind
 394 speed, cloud fraction) sampled at a given time can be used, such as the value at a given

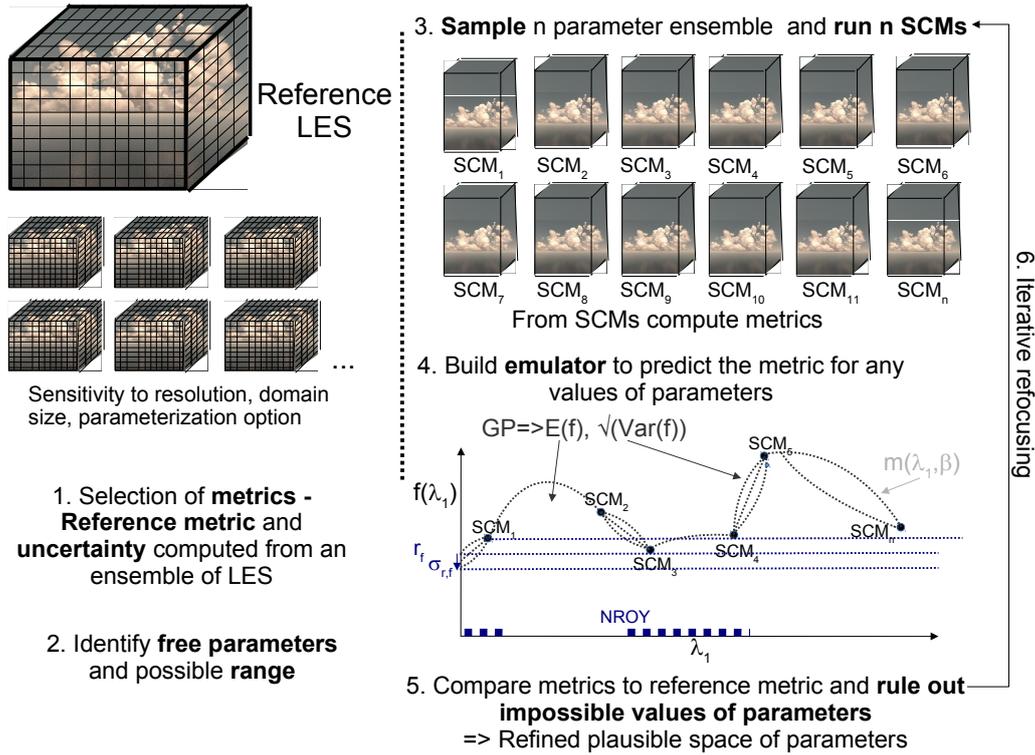


Figure 1. Schematic of the different steps of the htexplo tool

vertical level, the average or the maximum over a given layer (e.g., boundary layer, cloud layer), or the maximum over the whole atmospheric column. Radiation-oriented metrics are particularly relevant to enhance the link between the present process-oriented model calibration and the calibration of the corresponding 3D configuration. Ideally, the chosen metric should be as insensitive as possible to the model vertical resolution. In that regard, integrals (or averages) are good candidates for scalar metrics, as will be illustrated in Part II. Root-mean square errors are not encouraged for two reasons, i/ there are usually associated to a smaller signal to noise ratio and ii/ the implausibility (see section 3.6) is already a kind of root-mean square error. The number of metrics to be used is generally of the order of ten, but it can be many more.

More complex metrics such as vertical profiles, time series or spatial fields, can also be considered. In that case, methods are used to reduce the dimensions of the outputs and principal component decomposition is one option (e.g., Salter et al., 2019). However, scalar metrics, taken at a given time, or averaged over a short period of time, seem

409 often sufficient to robustly constrain most of the SCM simulations. Therefore, in the present
 410 paper and in Part II, only scalar metrics will be used.

411 References and their associated uncertainty are estimated from an LES ensemble.
 412 There are a priori two possibilities to build such an ensemble, which can be combined.
 413 The first consists in building the ensemble from simulations performed by different large-
 414 eddy models, as has been done in several GCSS intercomparison exercises (Brown et al.,
 415 2002; Siebesma et al., 2003; Stevens et al., 2005; vanZanten et al., 2011; de Roode et al.,
 416 2016). The reference thus corresponds to the LES ensemble mean, while the uncertainty
 417 is quantified by the LES ensemble variance. The second option, used in this paper, re-
 418 lies on only one large-eddy model and estimates the uncertainty around the reference model
 419 configuration by performing sensitivity experiments to horizontal and vertical resolution,
 420 domain size, and parameterization options (e.g., turbulence, microphysics, surface fluxes,
 421 radiation). In this study, we have chosen to use the simulation realized with the higher
 422 resolution over the largest domain and with the most relevant parameterization options
 423 as the reference, but the ensemble mean could also be used. The large-eddy model is the
 424 LES-configuration of Meso-NH (Lac et al., 2018). It makes use of a fourth-order centered
 425 discretization associated with an explicit fourth-order Runge-Kutta time integration. Fig-
 426 ure 2 illustrates the spread obtained from a Meso-NH LES ensemble exploring the sen-
 427 sitivity to horizontal, vertical resolution, domain size and options in the turbulence and
 428 cloud schemes for one given case, namely the ARM Cumulus case, which is a golden case
 429 for the study of continental cumulus (Brown et al., 2002). Table A2 in the Appendix de-
 430 scribes the different simulations used to estimate the uncertainty. Consistently with the
 431 literature (Brown et al., 2002; Matheou et al., 2011; vanZanten et al., 2011; Zhang et al.,
 432 2017), domain-average conserved thermodynamical quantities are weakly sensitive to changes
 433 in resolution, domain size and parameterization choices while the domain-average liq-
 434 uid water content and cloud fraction exhibit more spread. Metrics derived from those
 435 latter quantities will therefore be associated to a larger uncertainty. Figure 2 also indi-
 436 cates in grey shading the spread obtained from the LES intercomparison of Brown et al.
 437 (2002) highlighting a similar uncertainty estimate between the two methods mentioned
 438 above. Similar results are obtained for LES ensembles of other intercomparison exercises
 439 (not shown). For a given metric f , r_f is the reference metric value, estimated from the
 440 reference LES simulation or the average of the LES ensemble and $\sigma_{r,f}^2$ is the associated
 441 square error estimated from the LES ensemble. Note that, in the absence of available

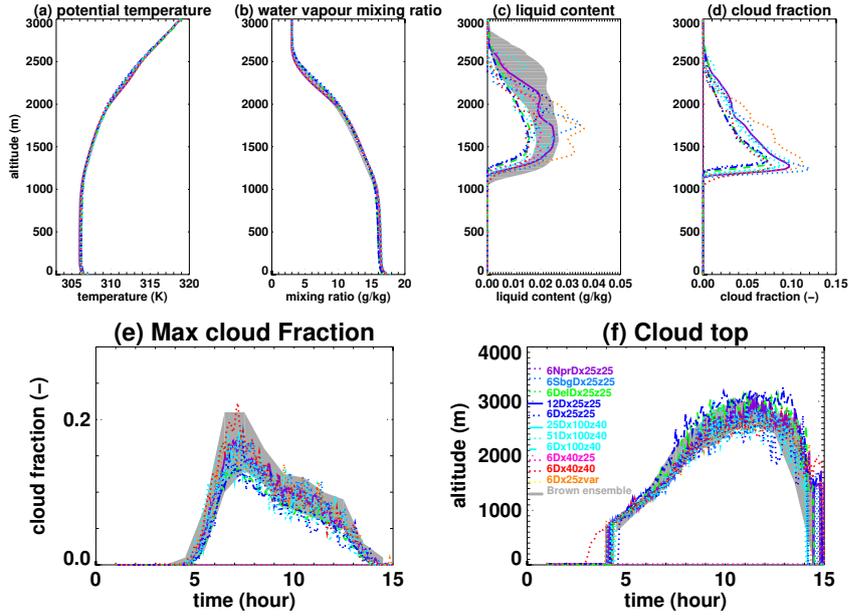


Figure 2. Vertical profile of (a) potential temperature, (b) water vapor mixing ratio, (c) liquid water content and (d) cloud fraction averaged over the horizontal domain at the 10th hour of the simulation (1530 LT) and time series of (f) the cloud top and (e) the maximum cloud fraction over the atmospheric column. The grey shading corresponds to the results of the Brown et al. (2002) intercomparison. The different color lines correspond to different sensitivity tests realized with Meso-NH changing either, one by one, the size of the domain, the vertical or horizontal resolution and some option in the cloud scheme, microphysics scheme or turbulence scheme (detailed in Table A2).

442 LES, observations can also be used as a reference to be compared to the SCM runs as
 443 illustrated in Ahmat Younous et al. (2018) but the observation error needs to be quan-
 444 tified.

445 3.3 Step 2: Selection of model parameters

446 The number of model parameters can be large (generally on the order of 10 for each
 447 parameterization). Estimating the prior range of values that needs to be explored for
 448 each of them requires the modeler's expertise. The definition of this range is an impor-
 449 tant step as the results are only valid in this predefined parameter space (Williamson

450 et al., 2013). So, we advise to choose a range as wide as possible in the absence of phys-
 451 ical reasons or numerical concerns for constraining it. Nevertheless, the user might con-
 452 sider some tradeoff as the smaller the ranges, the smaller the space to explore.

453 As the tool samples any parameter independently from the others (see Step 3), the
 454 method remains efficient even though a parameter with no influence on the results was
 455 included. A sensitivity analysis (Oakley & O’Hagan, 2004) could be used as a prelim-
 456 inary step in order to reduce the number of selected parameters but may not be a good
 457 idea in general (see section 3.8). The user can consider either linear or logarithmic vari-
 458 ations of the parameter values.

459 In the following, we consider a set of parameters $\boldsymbol{\lambda} = (\lambda_k)$, where the k param-
 460 eters are a subset of the model parameters involved in the different parameterizations
 461 (see section 1).

462 **3.4 Step 3: Experimental design and SCM runs**

463 Once the model parameters are selected and their range of values defined, an ex-
 464 perimental design is built. It corresponds to the selection of a relatively small set of val-
 465 ues for the model parameters $(\boldsymbol{\lambda}_i)_{i=1,\dots,n}$, usually on the order of ten times the number
 466 of parameters, as discussed in Loeppky et al. (2009). It explores the initial (or input)
 467 space of the parameter values in the range given for each parameter. An SCM simula-
 468 tion is performed for each of them and provides the state vector $\boldsymbol{x}_c(\boldsymbol{\lambda}_i)$. The objective
 469 is to "fill" the parameter space as uniformly as possible maximizing the minimum dis-
 470 tance between points. Here, as classically used for the design of computer experiments,
 471 a Latin Hypercube (LHC) (Williamson et al., 2015) is used to efficiently sample the in-
 472 put parameter space. Classically, a LHC for a n -member ensemble uniformly divides each
 473 dimension of the input space into n bins that are sampled once each and only once. All
 474 the parameters are thus varied simultaneously in contrast to other sensitivity analysis
 475 approaches such as in the Morris sensitivity analysis (Saltelli, 2002), where parameters
 476 are varied one by one. The LHC sampling used here maximizes the minimum distance
 477 between the selected points of the input space.

478 More precisely, here we use k -extended latin hypercubes as proposed by Williamson
 479 (2015). It consists in producing several LHCs, added sequentially, which ensure that each
 480 additional LHC samples an area of the space that has not been sampled yet by the pre-

481 various LHCs. Such a design provides the advantage of being able to robustly check the
 482 GP performance on well-designed sub-LHCs.

483 3.5 Step 4: Building emulators

484 The selected metric (see Step 1) is computed for each SCM simulation, noted $f(\boldsymbol{\lambda}_i)$
 485 for $i = 1, \dots, n$. These numbers serve as a training dataset for the building of an em-
 486 ulator. The emulator is then used to predict the metric values $f(\boldsymbol{\lambda})$ for any vector of pa-
 487 rameter values $\boldsymbol{\lambda}$ in the input space. A separate emulator is constructed for each met-
 488 ric.

Specifically, we use a Gaussian process (GP), a well known statistical model which has the advantage of interpolating observed model runs and provides a probabilistic prediction. The emulator gives a probability distribution for f written as:

$$f(\boldsymbol{\lambda}) \mid \boldsymbol{\beta}, \sigma^2, \boldsymbol{\delta} \sim \text{GP} (m(\boldsymbol{\lambda}, \boldsymbol{\beta}), k(\cdot, \cdot, \sigma^2, \boldsymbol{\delta})),$$

where $m(\boldsymbol{\lambda}, \boldsymbol{\beta})$ is a prior mean function with parameters $\boldsymbol{\beta} = (\beta_i)_i$ and k a specified kernel (a covariance function describing the covariance between any 2 points). The kernel has a parameter that normally controls variance, σ^2 , and parameters δ_k for each dimension of the input parameter λ_k that control the correlation attributed to each input. To start with, we assume a stationary kernel, i.e., the covariance only depends on the distance between points and not the absolute position. The GP is such that any finite collection $f(\boldsymbol{\lambda}_1), \dots, f(\boldsymbol{\lambda}_n)$ has a multivariate normal distribution with mean vector $m(\boldsymbol{\lambda}_1, \boldsymbol{\beta}), \dots, m(\boldsymbol{\lambda}_n, \boldsymbol{\beta})$, and variance matrix $\boldsymbol{\Sigma}$ with $\Sigma_{ij} = k(\boldsymbol{\lambda}_i, \boldsymbol{\lambda}_j, \sigma^2, \boldsymbol{\delta})$. Let the training data be $\mathbf{F} = (f(\boldsymbol{\lambda}_i))_{i=1, \dots, n}$, then

$$f(\boldsymbol{\lambda}) \mid \mathbf{F}, \boldsymbol{\beta}, \sigma^2, \boldsymbol{\delta} \sim \text{GP} (m^*(\boldsymbol{\lambda}, \boldsymbol{\beta}), k^*(\cdot, \cdot, \sigma^2, \boldsymbol{\delta})),$$

489 where there are well-known closed form expressions for m^* and k^* (Williamson et al.,
 490 2017). Note that m^* and k^* are the updated mean and covariance representing what the
 491 emulator has ‘learned’ from the data, \mathbf{F} .

492 Whilst there are many possible prior choices of m and k , htxplo uses a 2-phase
 493 approach. First, we impose a structured mean surface $m(\boldsymbol{\lambda}, \boldsymbol{\beta}) = \boldsymbol{\beta}^T \mathbf{g}(\boldsymbol{\lambda})$ as a linear
 494 combination of simple functions of the input parameters contained in the vector $\mathbf{g}(\boldsymbol{\lambda})$
 495 (e.g. monomials, Fourier functions and interaction terms are chosen through the forwards
 496 selection and backwards elimination method described in Williamson et al., 2013)). In

497 the second stage, we use the squared exponential kernel function and Hamiltonian Monte
 498 Carlo (HMC, implemented in Stan – Carpenter & Coauthors, 2017) to sample from the
 499 posterior distribution of the parameters $\boldsymbol{\beta}$, σ^2 , and $\boldsymbol{\delta}$ given \boldsymbol{F} (note that the mean sur-
 500 face $m(\boldsymbol{\lambda}, \boldsymbol{\beta})$ is not directly fitted in phase 1, but its structure is chosen, with Bayesian
 501 inference ultimately used in fitting for phase 2).

502 The choice of HMC implemented in Stan was motivated by requiring robust au-
 503 tomation of emulator building across many metrics and cases. Stan affords us with the
 504 ability to specify flexible and intuitive priors, and we use weakly informative priors as
 505 advocated by Gelman (2006). With the exception of the intercept term (which is uni-
 506 form), our prior for each $\boldsymbol{\beta}$ is $N(0, 10)$ and we use the ordinary least squares (OLS) fit-
 507 ted values as starting values for the HMC. We set $\delta_k \sim \text{Gamma}(4, 4)$ for all k to allow
 508 a wide range of potential correlation structures (this is a weakly informative prior) whilst
 509 penalizing very small values that typically have high likelihoods, but lead to emulators
 510 with no predictive power (for discussion, see Volodina, 2020). Our prior for σ^2 is a trun-
 511 cated Normal (at 0), with mean at the residual from our OLS fits, and variance set us-
 512 ing the variability of the ensemble (full details for these choices in Volodina, 2020).

513 The emulator is then tested using standardized Leave One Out diagnostics (e.g.
 514 Rougier et al., 2009) on the training data. These tests remove one point at a time from
 515 the training set and use the emulator fitted on the remaining data to predict the removed
 516 point. Repeated over the training set, we then check whether the majority of left out points
 517 lie within 95% prediction intervals (we would expect 5% to miss). Another check con-
 518 sists in removing a subdesign of the training set and attempting to predict it based on
 519 the new reduced training set. If the emulator fails these checks we revisit the computa-
 520 tion of the emulator. For example, the procedure described in Volodina and Williamson
 521 (2020) (and available in `htexplo`) can be used to derive an appropriate non-stationary
 522 kernel k before refitting the emulator by HMC. Once fitted, the GP expectation $E[f(\boldsymbol{\lambda})]$
 523 provides an estimation of the metric for any given $\boldsymbol{\lambda}$, and its variance $\text{Var}[f(\boldsymbol{\lambda})]$ provides
 524 an uncertainty around this estimation.

525 SCM runs are computationally cheap, but the fitted emulators are even cheaper
 526 and thus allow the computation of millions of predictions, with associated uncertainties,
 527 in a short time (a few minutes). This enables us to numerically define the space contain-
 528 ing acceptable sets of parameters with respect to the chosen metrics and in particular,

529 to visualize it (Step 5). The choice of Stan has proven effective for this project, though
 530 it does not scale well to larger ensembles. Going forward, a new version of the tools de-
 531 faulting to MAP estimation and using efficient parallel implementation has just been re-
 532 leased enabling millions of predictions in just a few seconds (Williamson & Volodina, 2020).

533 3.6 Step 5: History matching

The htexplo tool relies on the history matching technique, which seeks to rule out parameter values from the input space that are “implausible”, given the SCM behavior for these parameter values and the sources of uncertainty. These sources include the reference (observation) error, treated as a random quantity with mean 0 and variance $\sigma_{r,f}^2$, and the SCM discrepancy, which has mean 0 (unless the user knows the direction in which the model is biased) and variance $\sigma_{d,f}^2$ (Sexton et al., 2011). The emulator is used to estimate the model behavior on a much larger sample of the input space than possible with the SCM. To history match the SCM behavior, we introduce the “implausibility” measure for the metric f (Williamson et al., 2013), $I_f(\boldsymbol{\lambda})$, which is a distance between the metric prediction $f(\boldsymbol{\lambda})$ by the emulator at $\boldsymbol{\lambda}$, and the reference metric value, r_f , with respect to the norm induced by our second-order uncertainty specification, noted $\|\cdot\|_H$ below. The implausibility reads

$$\begin{aligned}
 I_f(\boldsymbol{\lambda}) = \|r_f - f(\boldsymbol{\lambda})\|_H &= \frac{|r_f - \mathbf{E}[f(\boldsymbol{\lambda})]|}{\sqrt{\text{Var}[r_f - \mathbf{E}[f(\boldsymbol{\lambda})]]}} \\
 &= \frac{|r_f - \mathbf{E}[f(\boldsymbol{\lambda})]|}{\sqrt{\sigma_{r,f}^2 + \sigma_{d,f}^2 + \text{Var}[f(\boldsymbol{\lambda})]}}.
 \end{aligned}
 \tag{5}$$

534 The model discrepancy for the metric f , $\sigma_{d,f}$, accounts for the model structural
 535 error due to the inherent inability of the SCM to reproduce the LES exactly (due to un-
 536 resolved physics or missing processes, for example). It could be defined as the minimum
 537 error possible when exploring the full set of parameters, however, this could permit the
 538 SCM to be close to the reference for the wrong reasons and does not account for mul-
 539 tiple metrics and cases, so we avoid this definition. Instead it is typically defined to be
 540 the uncertainty left in the difference between the SCM metric when the parameters are
 541 fixed at their best values (fixed the same for all metrics) and the references. This quan-
 542 tity is perhaps the target of model development in the first place and, as such, is unknown.
 543 For example, suppose we want to test the ability of a new parameterization to capture
 544 the behavior of the reference. With the standard definition of discrepancy, the uncer-
 545 tainty needed so that the new parameterization captures the behavior of the reference,

546 it is not clear how to proceed with testing. Our approach instead is to treat model dis-
 547 crepancy as a “tolerance to error” as detailed in Williamson et al. (2017). The tolerance
 548 to error is the distance between model results and the reference that the modeler would
 549 be satisfied with, enabling modelers to place confidence in certain metrics/parts of their
 550 parameterization, and relax restrictions on others as needed. As illustrated in section
 551 4 and Part II, defining this tolerance to error can be a difficult a-priori task; however ex-
 552 perimenting with this value provides important insights into the behavior of the model
 553 and its inherent limitations. The most attractive feature of this approach to discrepancy
 554 is that, for a given tolerance to error, if the induced NROY space is empty it means that
 555 the parameterization is not able to reproduce the reference under the given tolerance.
 556 Either the tolerance can be relaxed, accepting the limitations of the current set of pa-
 557 rameterizations, or the parameterization can be revisited.

The implausibility defines a membership rule for NROY space after the first iter-
 ation:

$$\text{NROY}_f^1 = \{\boldsymbol{\lambda} \mid I_f(\boldsymbol{\lambda}) < T\}.$$

558 where T is a chosen threshold (or cutoff). For scalar metrics, it is standard to use $T =$
 559 3 justified using Pukelsheim’s rule that states 95% of the probability density for any uni-
 560 modal distribution is within 3 standard deviations of the mean (Pukelsheim, 1994). Us-
 561 ing this threshold makes it unlikely that good parameter values are ruled out by chance.
 562 To measure and visualize NROY space the implausibility $I_f(\boldsymbol{\lambda})$ is calculated on a ran-
 563 dom LHC sampling of a large number (on the order of hundreds of thousands or millions)
 564 of vectors $\boldsymbol{\lambda}$.

565 Note that $I_f(\boldsymbol{\lambda})$ can be smaller than the chosen threshold T either because $E[f(\boldsymbol{\lambda})]$
 566 is close to the reference or because the sum of the different errors is large. When the un-
 567 certainty of the emulator is larger than the tolerance to error and observation error, points
 568 that should be ruled out are kept in the NROY. In this case, further iterations are de-
 569 sirable in order to increase the density of the sampling of NROY and hence improve the
 570 emulator quality and reduce the associated uncertainty.

571 **3.7 Iterative refocusing and multi-metrics**

572 One advantage of this method is to progressively optimize the design of simulations
 573 to be run. New simulations are iteratively added only where it is useful to increase the

574 emulator accuracy. This is performed by iterating the same process previously described
 575 several times in "waves", (this is termed "iterative refocusing" and is a fundamental part
 576 of the history matching approach). Each new iteration n starts from the remaining space
 577 NROY_f^{n-1} estimated at the end of the previous wave. Because of its complex geometry,
 578 a LHC sampling, as in the first wave, cannot be applied, and therefore the remaining space
 579 is re-sampled uniformly. A new SCM simulation ensemble is performed with this design
 580 and is used to proceed with steps 4 and 5. The new emulator is only valid in the new
 581 parameter space, namely NROY_f^{n-1} . Outside this space, we rely on the emulators from
 582 the previous waves. As in Step 5, to measure and visualize NROY_f^n , the implausibility
 583 is computed over a large number of points in the input space. The threshold T may be
 584 varied between waves, but we advise to keep it to 3 as long as the process has not con-
 585 verged (i.e. the emulator variance within the current NROY space remains large – see
 586 also section 4 and Part II). The iterative refocusing stops when the convergence of the
 587 sequence $(\text{NROY}_f^n)_n$ has been qualitatively achieved.

So far, we have considered only one metric, but several metrics $(f_k)_k$ can be com-
 bined at the same time. An implausibility is then computed for each metric and the to-
 tal NROY^n space is the intersection of the $\text{NROY}_{f_k}^n$ associated with each metric:

$$\text{NROY}^n = \bigcap_k \text{NROY}_{f_k}^n = \{ \boldsymbol{\lambda} \mid \#\{k \mid I_{f_k}^n(\boldsymbol{\lambda}) > T\} \leq \tau \},$$

588 $\#$ represents the number of metrics fulfilling the condition indicated into brackets (where
 589 the implausibility is greater than the threshold) and τ , the number of metrics for which
 590 the model is allowed to be far from the reference while still kept in the NROY space. If
 591 $\tau = 0$, all metrics must satisfy our implausibility cutoff. If there are a large number of
 592 metrics then τ should be increased ($\tau \geq 1$) to avoid multiple testing problems mean-
 593 ing that too many good parameter values are ruled out by chance. If a modeler seeks
 594 to prioritize certain metrics, they can either be introduced in early waves, ensuring that
 595 the NROY space satisfies priority metrics first before introducing new ones, or the tol-
 596 erance to error, which is defined for each metric, can be used to impose priorities (a larger
 597 tolerance to error induces a less constraining metric).

598 **3.8 Sensitivity analysis provided by the tool**

599 The htexplo tool provides its own sensitivity analysis, which, due to the use of multi-
 600 wave history matching, is rather different from traditional methods applied to models

601 throughout the literature (Bastidas et al., 2006; Guo et al., 2014; Johnson et al., 2015).
602 Traditional methods, either derivative-based (Saltelli, 2002), or variation-based (Oakley
603 & O’Hagan, 2004), essentially seek to identify which parameters modify model output.
604 This can help focus further study, model development or even observation collection to
605 help understand these parameters. Note that the htexplo tool provides at the first it-
606 eration a sensitivity analysis over the entire space where correlation among parameters
607 is included as the parameters are not varied one at a time.

608 However, for calibration purposes, once history matching is considered as a valid
609 approach for a given model, the sensitivity analysis should not be done on the full model
610 input space. By using history matching, we acknowledge that there is a large part of the
611 model parameter space that is not useful for understanding reality. The Gaussian pro-
612 cesses remove this uninformative space in order to target the space where the model be-
613 comes useful. Once we have this useful subspace, the usual and important questions that
614 are posed by sensitivity analysis should be considered. For example, how is the model
615 output changing as we move through parameter space and which parameters are respon-
616 sible for these changes? As will be illustrated in section 4, the NROY visualization al-
617 lows us to see, as we move in two dimensions of a parameter space, in addition to the
618 possible values of each parameters, which combinations of parameters it is important to
619 get right. As all models within the NROY space are consistent with our metrics, sen-
620 sitivity analysis as described here is now really focused on the relevant subspace. Note
621 that sensitivity analysis on the original input space does not answer these questions. Seen
622 through the history matching lens, on the full space, sensitivity analysis is showing us
623 which parameters are responsible for the variability in the space we are about to cut. Whilst
624 informative for helping us cut the space efficiently, sensitivity analysis is not necessary
625 at this stage. Our methods are already efficiently able to do this. As well as all of the
626 benefits we have for tuning, we would argue that history matching is achieving many of
627 the same things that a sensitivity analysis achieves in terms of informing the modeling,
628 but concentrated only on the model input space that is consistent with the observations.

629 Performing variance-based sensitivity analysis in NROY space is not trivial and we
630 are not aware of any methods that are currently able to do this. Variance-based sensi-
631 tivity analysis requires independent input spaces (which is what we always start with
632 in Wave 1). But after cutting space, we have complex relationships between the param-
633 eters. NROY space may not even be simply connected, and can be highly non-linear. Ef-

634 efficient methods for calculating sensitivity in these unusual spaces would be interesting
 635 to apply for history matching as an avenue for further research.

636 **3.9 On the use of history matching and the avoidance of optimization**

637 Whilst history matching is well established and is being used in a growing num-
 638 ber of climate studies, other methods of calibration are more popular and we believe should
 639 be avoided for process-based model development. Whilst many methods based on op-
 640 timizing a cost function exist (Hourdin et al., 2017), the most popular in the UQ com-
 641 munity is Bayesian calibration (Kennedy & O’Hagan, 2001). Bayesian calibration requires
 642 a similar set up to history matching (emulators, observation errors and model discrep-
 643 ancy) and then jointly finds the posterior probability distribution of the “best” value of
 644 the input parameters and the model discrepancy (strong prior information on the dis-
 645 crepancy is required to make this sensible, Brynjarsdóttir & O’Hagan, 2014). Optimiza-
 646 tion methods like these do not afford us with the chance to falsify a parameterization
 647 (they always find the best value), nor do they give all parameter values that are consis-
 648 tent with the observations (in our case reference LES) that can then be used when tun-
 649 ing the 3D model (see Part II).

650 **4 Illustration of htexplo on a simple case**

651 In this section, the use of htexplo is illustrated for the ARPEGE-Climat 6.3 atmo-
 652 spheric model (Voldoire et al., 2019; Roehrig et al., 2020) based on a single 1D case. More
 653 comprehensive exploitation of the tool will be given in Part II.

654 **4.1 Model, parameters and case-study**

655 We use the SCM version of ARPEGE-Climat 6.3, the atmospheric component of
 656 the CNRM-CM6-1 climate model (Voldoire et al., 2019; Roehrig et al., 2020) and aim
 657 at analyzing the importance of the values of free parameters of the turbulence scheme
 658 (based on Cuxart et al., 2000) on the simulation of an idealized clear boundary layer.
 659 Details on the ARPEGE-Climat atmospheric component, the turbulence scheme and the
 660 used configuration are given in the Appendix B. Among the different free parameters of
 661 the turbulence scheme, three are selected for this analysis. A_ϵ controls the expression
 662 of the dissipation length-scale as a function of the mixing length-scale; A_U and A_T re-

663 spectively enter into the expression of the exchange coefficient in Eq. B3 for the wind
 664 and the temperature (the same coefficient, A_U , is used for both the zonal and meridional
 665 component of the wind). The range of variation explored for each parameter is indicated
 666 in Table 1 and the parameters are varied linearly in those ranges. The turbulence pa-
 667 rameterization includes other free parameters but the three most influential parameters
 668 for this case have been selected and no free parameters of the mass-flux scheme are con-
 669 sidered.

Table 1. List of the free parameters of the turbulence scheme that are varied in this example with default values and range of variation

Names	Default	Minimum	Maximum	Parameter Description
A_U	0.126	0.01	0.4	Affects the eddy-diffusivity of momentum
A_ϵ	0.85	0.1	3.	Controls the dissipation length-scale
A_T	0.14	0.01	1.	Affects the eddy-diffusivity of temperature

670 To keep the example simple, only one case is used here. This case is a dry ideal-
 671 ized case of a convective boundary layer with a constant-in-time large surface sensible
 672 heat flux of 270 W m^{-2} ($Q_* = 0.24 \text{ K m s}^{-1}$ in, Ayotte et al., 1996) with a strongly capped
 673 boundary layer, called 24SC in the following. The importance of combining different cases
 674 will be illustrated in part II.

675 We first document a sequence of three waves where additional metrics are added
 676 at each iteration (Experiment 1). We will then discuss the results obtained when adding
 677 all the metrics directly at Wave 1 (Experiment 2), varying the threshold used to deter-
 678 mine the NROY (Experiment 3 see also section 3.5), using more SCM runs (Experiment
 679 4), and varying the tolerance to error (Experiments 5 and 6).

680 4.2 Three consecutive waves adding metrics progressively

681 For the first iteration (or wave in the following) of Experiment 1, 30 SCM simu-
 682 lations of the 24SC case were realized by varying values for the three parameters explor-
 683 ing at best (using a LHC sampling, see section 3.4) the range of each parameters (Ta-
 684 ble 1). Figure 3 illustrates that the parameters are randomly sampled as indicated by

685 the distribution of the black dots along the different x-axes. Three different metrics are
686 used to characterize the turbulent mixing in the boundary layer and are progressively
687 introduced through the successive waves. The first chosen metric is the potential tem-
688 perature averaged over the layer 400-600 m. It is a good proxy for the boundary-layer
689 potential temperature, which is well mixed between the surface and the boundary-layer
690 top, located around 1300 m. This metric is computed for the 30 SCM runs; these com-
691 putations serve as training data for the construction of the emulator. The prior mean
692 function (see section 3.5), m , for this emulator is a sum of linear and quadratic functions
693 of the parameters. The stationary squared-exponential kernel provides a sufficient fit to
694 the data according to the leave-one-out methodology. Figure 3 presents the variation of
695 the metric as a function of the parameters: some first-order relationships appear with
696 the boundary-layer potential temperature increasing with A_U and A_T to a lesser extent
697 (due to an increased mixing associated to a larger diffusivity and larger fluxes) and de-
698 creasing with A_ϵ (due to a reduced mixing because of the increased dissipation). For this
699 metric, we have chosen a tolerance to error of 0.5 K. This may be a bit large for this very
700 idealized case (with no moisture, an already convective initial state) but this is an er-
701 ror we will be satisfied with generally for boundary-layer potential temperature. Given
702 this tolerance to error (indicated by the dashed horizontal grey line), the metric does not
703 provide much constraint on the model behavior and the entire initial parameter space
704 is kept (c.f. Table 2). Note that this tolerance to error is much larger than the uncer-
705 tainty around the LES ($\sigma_{r,f} = 0.075$ K) and the emulator (this uncertainty varies across
706 the values of the parameters; it is quantified here as the mean of the standard deviation
707 for all the points of the dataset during the LOO experiment. For wave 1 and the first
708 metric, it is 0.042 K). Section 4.3 details the effect of a reduced tolerance to error.

A second wave is realized, with 30 runs sampling the NROY space of the first wave (the previous 30 SCM runs could also have been used for efficiency), which is in fact the entire initial parameter space as the first metric did not constrain the parameter space. Two metrics are computed from those 30 runs: the potential temperature averaged between 400 m and 600 m as in the first wave and the entrainment metric, A , quantifying the overshoot of the boundary layer relative to the initial profile as defined in Ayotte et al. (1996). A is computed as:

$$A = \frac{\int_{z_i(t_0)}^H (\theta(z, t_f) - \theta(z, t_0)) dz}{t_f - t_0} = \frac{\int_0^H (\max(\theta(z, t_f) - \theta(z, t_0), 0)) dz}{t_f - t_0}$$

Table 2. Description of the model discrepancy (Disc.) of the given metric (indicated in the 2nd, 3rd and 4th columns), the Cutoff, threshold used for implausibility (5th column) and the Not-Ruled-out-Yet Space (fraction in % of initial space of parameters, 6th column) for each metric (7th column) for each Experiment and wave.

N^o Expt N^o Wave	$\sigma_{d,\theta_{BL}}$ [K]	$\sigma_{d,Ay\theta}$ [Kms ⁻¹]	$\sigma_{d,ws_{BL}}$ [m s ⁻¹]	Cutoff	NROY (%)
Exp1-1	0.5	-	-	3	100
Exp1-2	0.5	0.05	-	3	30
Exp1-3	0.5	0.05	1	3	23
Exp1-4	0.5	0.05	1	3	20
Exp1-5	0.5	0.05	1	3	18
Exp2-1	0.5	0.05	1	3	40
Exp2-2	0.5	0.05	1	3	38
Exp2-3	0.5	0.05	1	3	27
Exp2-4	0.5	0.05	1	3	17
Exp3-1	0.5	0.05	1	3	72
Exp3-2	0.5	0.05	1	3	32
Exp3-3	0.5	0.05	1	2.5	22
Exp3-4	0.5	0.05	1	2.	15
Exp4-1	0.5	0.05	1	3	25
Exp4-2	0.5	0.05	1	3	19
Exp5-1	0.25	0.025	0.5	3	32
Exp6-1	0.1	0.01	0.25	3	31

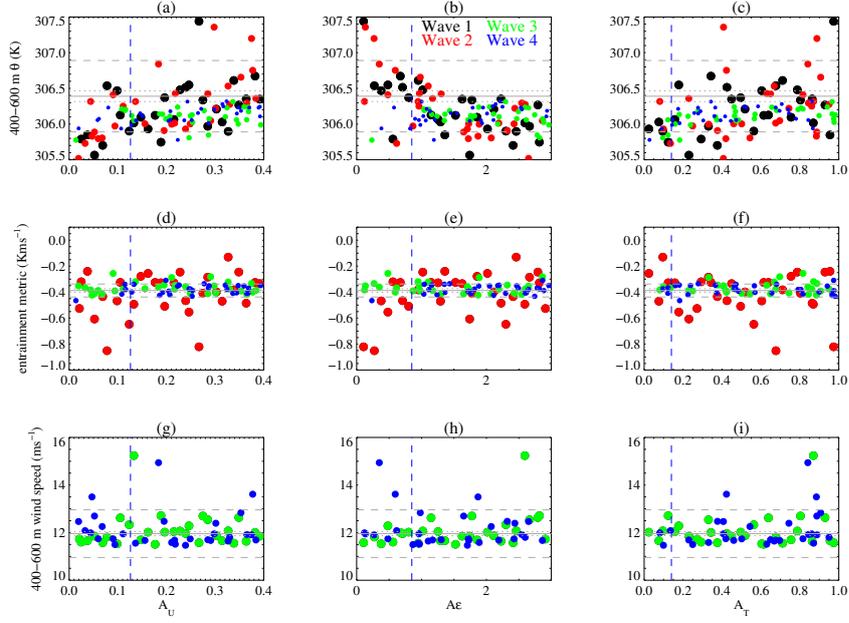


Figure 3. The three metrics, boundary-layer potential temperature (a–c), entrainment metric (d–f) and boundary-layer windspeed (g–i) are plotted as a function of the value of each parameter, A_U (a, d, g), A_ϵ (b, e, h) and A_T (c, f, i). A different color is used for the different waves of Experiment 1 (black for Wave 1, red for Wave 2, green for Wave 3 and blue for Wave 4). The vertical dashed blue line corresponds to the default value of the parameter used in the model, the horizontal thin full grey line correspond to the reference metric and the dotted lines indicates the uncertainty around this reference from the different LES simulations while the dashed lines indicate the tolerance to error around the reference.

709 t_0 being the initial time, t_f the time at which the metric is computed and H the top of
 710 the model or a level largely above the boundary-layer top. This metric is less commonly
 711 used for evaluating models and it was more difficult to specify a tolerance to error, which
 712 was taken as 0.05 K.m s^{-1} . An emulator is built for each metric. The second metric is
 713 more restrictive and the NROY space is now reduced to 30% of the initial parameter space
 714 (Table 2). The obtained NROY (not shown) is not very different from the one obtained
 715 for the third wave. It excludes values of the parameters that lead to simulations with
 716 too large or too small entrainment metric as indicated by the differences between the red
 717 dots and the green ones in Fig. 3.

718 A third wave is realized, with 30 new SCM runs sampling the new NROY. Three
 719 metrics are computed from those 30 runs: the two previous ones plus the wind speed av-
 720 eraged between 400 m and 600 m. For this last metric, we fixed the tolerance to error
 721 to 1 m s^{-1} . After this third iteration, the NROY is 23% of the initial space. As shown
 722 in Fig. 4, the spread of the different simulations that sampled the parameter values re-
 723 duces progressively throughout the different waves and this tool allows to discard val-
 724 ues of parameters that induce a too deep boundary layer. The wind-speed profiles did
 725 not completely converge and this is associated to the tolerance to error, which has been
 726 fixed to 1 m s^{-1} .

727 The uncertainty around the LES obtained from eight different LES runs with slightly
 728 different configurations, detailed in the appendix A, is 0.075 K for θ_{BL} , 0.014 K m s^{-1}
 729 for A_θ and 0.083 m s^{-1} for ws_{BL} , on the same order of magnitude of the emulator un-
 730 certainty. For the first and third metrics, the tolerance to error is much larger than the
 731 reference and emulator uncertainties while for the second metric the three uncertainties
 732 are of the same order of magnitude.

733 The final NROY space after the third wave is visualized in Fig. 5. This figure shows,
 734 on the upper right side, the two-dimensional density plots of the acceptable parameter
 735 space for each pair of parameters. For a given point in each sub-figure the shading in-
 736 dicates the percentage of the domain in the other dimensions ($n-2$, here only one as only
 737 three parameters are considered) that is acceptable. The metrics tend to reject prefer-
 738 entially low values of A_ϵ with high values of A_U or high values of A_ϵ with low values of
 739 A_U underlying some correlation between these two parameters. As a practical tool, those
 740 density plots provide their own type of second-order sensitivity analysis. They allow us

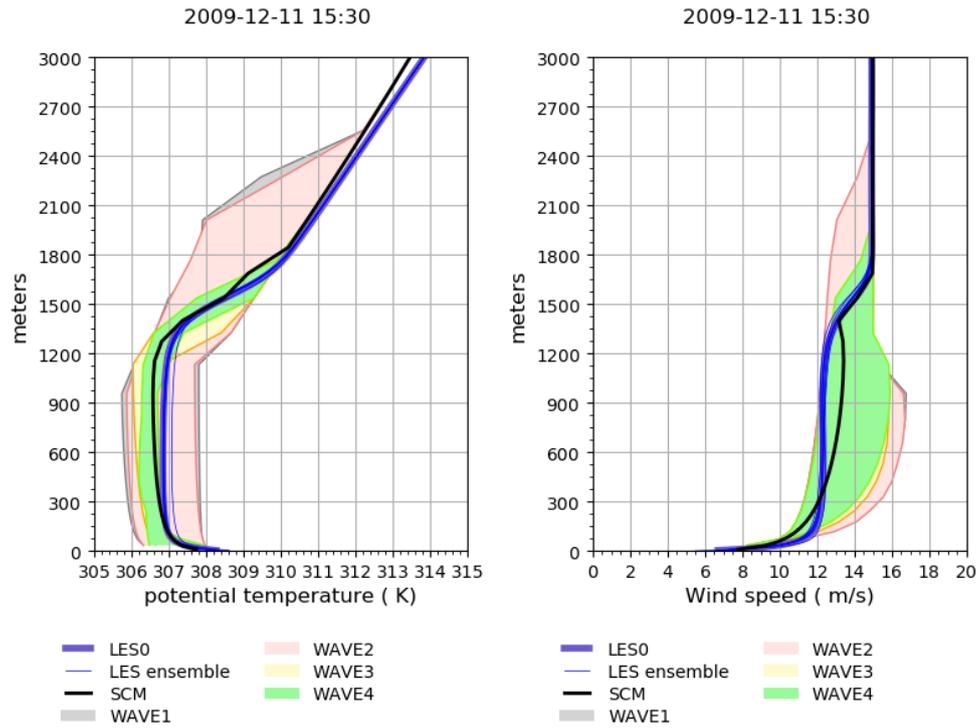


Figure 4. Vertical profile of (a) potential temperature and (b) wind speed for the last hour of the simulation with the spread of the ensemble of simulations used for the different waves indicated in different color shadings for Exp 1, the default simulation is in black, the reference LES in thick dark blue and the different elements of the LES ensemble in thin blue lines.

741 to see, as we move in two dimensions of the parameter space, how the shape is chang-
 742 ing and, moreover, which combinations of parameters it is important to get right and,
 743 not usually included in a sensitivity analysis, how they need to be set in order to get sen-
 744 sible answers. The default values of the parameters are within the NROY space confirm-
 745 ing that they correspond to an acceptable calibration of the turbulence scheme, given
 746 the chosen tolerance to error and the LES uncertainty. This is also confirmed by the sim-
 747 ulations of the last wave having a behavior similar to the default simulation as shown
 748 in Fig. 4.

749 4.3 Robustness

750 In this subsection, we analyze the sensitivity of the results to i) the sequence of in-
 751 troduction of metrics (Experiment 2 uses the three metrics directly at Wave 1), ii) the
 752 threshold used to determine the NROY space (Experiment 3), iii) the number of SCM
 753 runs used to form the training dataset (Experiment 4), and, iv) the tolerance to error
 754 (Experiments 5 and 6).

755 If the three metrics are introduced directly in the first wave (Experiment 2), the
 756 NROY space is similar in shape to the one obtained after three waves (see Table 2 and
 757 Fig. 5) although the NROY space is larger (40% against 23%). Repeating more waves
 758 with the same metrics allows to progressively converge to the same NROY space. Note
 759 that a test with only one metric but the most constraining one, namely the entrainment
 760 metric, leads to very similar result ($NROY = 43\%$) for the first wave (not shown). Al-
 761 though not illustrated for this case, introducing the metrics one by one, is sometimes im-
 762 portant: i/ it can allow us to give some priority among the metrics, first finding a space
 763 consistent with the first metric in which the second metric is then used as a constraint
 764 and ii/ if one metric has a strong non-linear behavior reducing the initial parameter spaces
 765 with other metrics may increase the capacity of the emulator to reproduce the metric
 766 behavior. These results also indicate that adding a new metric in the core of the pro-
 767 cess does not alter the selection, allowing us to add supplementary metrics if one real-
 768 izes that some behavior of the SCM is not constrained enough, a fundamental aspect of
 769 history matching. Defining when to stop the iteration is not easy. We recommend to stop
 770 iterations when the NROY stops to significantly decrease. At this stage, one can reduce
 771 the cutoff used to define the implausibility and re-iterate with this new cutoff. This is
 772 illustrated with more detail in part 2. Here, Table 2 shows that a NROY of 18% is ob-

773 tained after Wave 5 for Experiment 1, Wave 4 for Experiment 2 or Wave 2 for Exper-
 774 iment 4. We can assume that for this cutoff the convergence is reached at those waves.

775 In Experiment 3, we first realize two waves as in Experiment 2 and then progres-
 776 sively reduce the threshold used to determine the NROY space from 3 to 2.5 in Wave
 777 3 and from 2.5 to 2 in Wave 4 (see Table 2) to explore the impact of less conservative
 778 threshold (a threshold of 3 corresponds to ruling out what exceeds three times the un-
 779 certainties and keeps 95% of the probability for any unimodal probability distribution).
 780 The differences in the NROY space of the first wave with Exp2-1 indicates that 30 SCM
 781 runs are probably not enough to robustly constrain the first iteration and more itera-
 782 tions are needed. Then, reducing the cutoff induces a smaller NROY space but the change
 783 is not radical. This was expected from the lower left figures of Fig. 5 that show the min-
 784 imum value of the implausibility for any variations of the other parameters (here, the
 785 third parameter). Indeed, the area with minimum value of $I_f(\boldsymbol{\lambda}) > 3$ (i.e. the points
 786 that are excluded from the NROY space whatever the value of the third parameter) is
 787 very similar to the area with minimum value of $I_f(\boldsymbol{\lambda}) > 2$.

788 All of the previous experiments have been realized using a rather small training dataset
 789 of 30 SCM runs (ten times the number of parameters). Experiment 4 has tested the im-
 790 pact of using 90 SCM runs instead of 30 for wave 1. This experiment produces directly
 791 a smaller NROY space (NROY=25%; Fig. 6) at the first wave than obtained from 30 SCM
 792 runs (see Exp3-1 or Exp2-1 in Table 2). A compromise must be found between a larger
 793 ensemble of simulations that increases robustness but is more costly.

794 The sensitivity to the tolerance to error is illustrated in Table 2 and Fig. 6 with
 795 Experiments 5 and 6. When reducing the tolerance to error by a factor of two the NROY
 796 space is 32% of the initial space in Exp5-1 (using the three metrics at once, so to be com-
 797 pared to 40%). The NROY space (31% of the initial space) is not much reduced further
 798 when reducing the tolerance to error twice more (Exp6-1), because the tolerance to er-
 799 ror is not anymore the limiting uncertainty. It is interesting to note that even when strongly
 800 reducing the tolerance to error, the default values for the three selected parameters are
 801 still in the NROY space validating the choice of parameter values used in the control sim-
 802 ulation. The lower left panel of the subfigures in Fig. 5 and Fig. 6 indicates the mini-
 803 mum implausibility along the other dimensions of the space and as illustrated in Fig. 6,

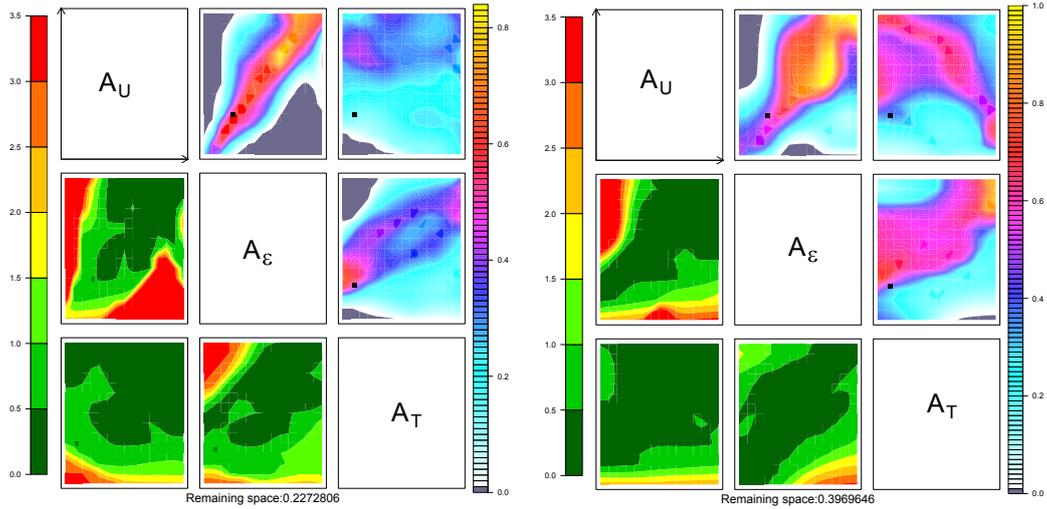


Figure 5. The left panel corresponds to the result of Exp1-3 and the right panel to Exp2-1. The upper right triangle contains 3 subfigures showing 2D sub-matrix. Each sub-matrix is a restriction to 2 parameters, the name of which are given in the diagonal of the main figure, and presents in colors the fraction of points with implausibility smaller than the threshold (here a value of 3). This fraction is obtained by fixing the two parameters at values of the x-axis and y-axis of the plotted location and searching the other dimensions (here the third dimension as we have only three parameters) of the parameter space. This allows to visualize in 2-D the full NROY which is 3-D here but can be n-D if n parameters are selected. The lower left triangle (with also 3 subfigures) presents the minimum value of the implausibility when all the parameters (here only one) are varied except those used as x- and y-axis. These plots are orientated the same way as those on the upper triangle, for easier visual comparison. The black dots correspond to the default values used in the model.

804 reducing the tolerance error (when larger than the other errors) induces a reduction of
 805 the denominator in the implausibility and therefore an increase of implausibility.

806 **5 Conclusion**

807 In this paper, we make a proposal to accelerate weather and climate model devel-
 808 opment. Our proposal tackles model development and calibration jointly. For that pur-
 809 pose, we have developed a tool that formalizes a process-based calibration, the *High-Tune*
 810 *Explorer* made available to the other modeling groups. It extensively exploits the SCM/LES
 811 comparison on a multicases, multi-metrics basis and benefits from machine learning tech-

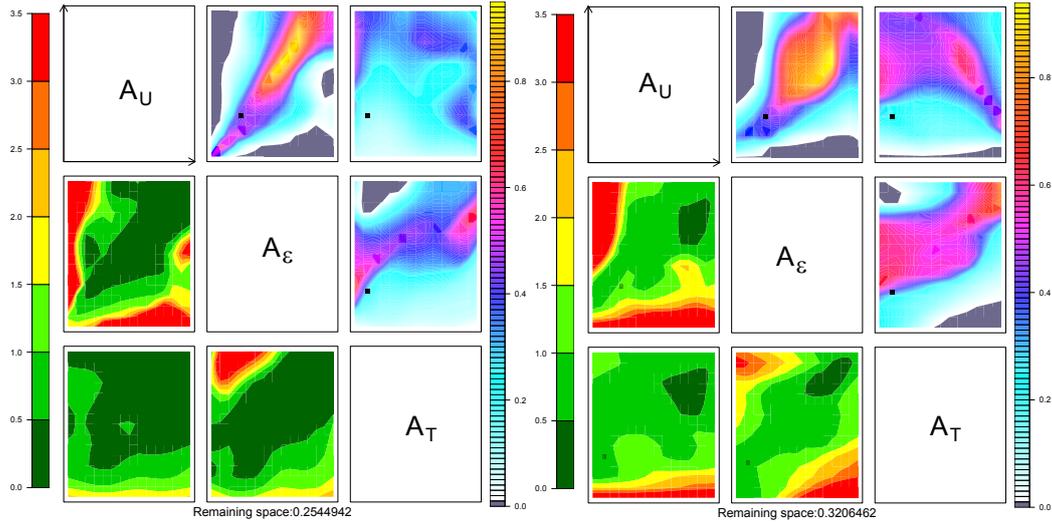


Figure 6. Same as Fig. 5 but for the sensitivity to the number of SCM runs (Experiment 4, left panel) and to the tolerance error (Experiment 5, right panel).

812 techniques. In contrast with other recent proposals to use machine learning techniques in cli-
 813 mate modeling, we keep parameterizations as key ingredients of these models because
 814 they summarize our current understanding of the main physical processes. This choice
 815 is motivated in particular by the confidence needed when extrapolating the model re-
 816 sults to a future climate.

817 The tool allows us to define the sub-domain of the parameter values for which SCM
 818 matches LES on selected metrics for a series of cases within a given uncertainty. The ex-
 819 ploration of the free-parameter space is facilitated using Gaussian process emulators. These
 820 emulators, once trained on a limited number of real simulations, predict the SCM with
 821 uncertainty for any value of the parameters in a much shorter time than required to run
 822 the SCM. History matching using the emulator is performed iteratively to progressively
 823 shrink the space of acceptable parameter values. This iterative approach contrasts with
 824 the more traditional tuning strategy based on optimization, which seeks an individual
 825 “best” value where the SCM minimizes a cost function computed for a set of given met-
 826 rics. The latter approach strongly depends on the weights given to each metric and is
 827 highly sensitive to the choice of metrics. By pursuing a strategy for discarding param-
 828 eter values, we are left with a free parameter domain that is (i) consistent with the met-
 829 rics we have chosen, (ii) can be further reduced by introducing new metrics or altering

830 our tolerance to model error, and (iii) does not claim a single best simulation which may
831 be over-fitted to one or more metrics, needlessly biasing the simulation and potentially
832 leading to less physical behavior than other choices in our not-ruled-out-yet space when
833 the model is projected into different regimes. Our tool formalizes the consideration of
834 the different sources of uncertainties associated to the reference, the statistical tool and
835 the model. For the latter, we take a “tolerance to error” approach, allowing the ques-
836 tion of whether a parameterization can match our reference as well as we think it ought
837 to, and enabling us to understand the model’s limitations throughout the process.

838 In the present study, we present applications of the *High-Tune Explorer* to the SCM/LES
839 framework, focused on the representation of the atmospheric boundary layer. We have
840 illustrated how this tool allows us to objectively verify choices that have been made by
841 model developers for the free-parameter values. Experimenting with the combination of
842 the metrics with this tool allows us to clarify the importance of a given metric, the num-
843 ber or combination of metrics that should be used, and the possible redundancy between
844 metrics all in an efficient way that was not possible before. The tool also enables us to
845 include new metrics at a new iteration so that we can pursue the calibration exercise,
846 even though one realizes an important deficiency of the model is not addressed by the
847 previously selected metrics. Our framework allows a progressive addition of metrics, cases
848 or a gradual reduction of the tolerance to error and is therefore very flexible.

849 Although this new framework is tested here for the improvement of boundary-layer
850 processes (turbulent transport in Part I and cloud representation in Part II) by running
851 the full atmospheric physics on one model column considering well established test cases
852 for which LES are particularly relevant, it has much broader application. It can be used
853 for instance to calibrate elementary pieces of parameterization (e.g., entrainment formu-
854 lation) without time integration. This methodology can be easily expanded to other pa-
855 rameterizations as well. The key ingredient for doing this is a reliable reference with doc-
856 umented uncertainty. This reference could come either from a detailed modeling of the
857 process, as done here with LES, or from observations as long as the other sources of dis-
858 crepancy, as the uncertainty coming from the case definition, are documented. Propos-
859 ing new relevant metrics and estimation of associated uncertainties will become valuable
860 now that we know how to include them in the model improvement process. An effort is
861 currently done in that direction in parallel to the work presented here, consisting in pro-
862 viding reference radiative transfer computations on the classical cloud test cases currently

863 used for parameterization development. The development of the parameterization of bound-
864 ary layer and clouds based on SCM/LES comparisons focused so far on the representa-
865 tion of atmospheric transport and macrophysics of clouds, but the radiative transfer com-
866 putations run in LES models were often not more reliable than those used in GCM, pre-
867 venting the use of radiative metrics. By developing fast and accurate radiative tools that
868 account for the full 3D radiative transfer in LES cloud scene, as proposed by Villefranque
869 et al. (2019), we can compute many types of radiative metrics, from monochromatic, lo-
870 cal, and directional observable to integrated energetic quantities. The use of such radia-
871 tive metrics will allow us to tackle calibration of radiative parameterizations but also to
872 better link the calibration realized at the level of the parameterizations itself with the
873 one realized for the final full 3D model calibration, which mainly targets the radiative
874 forcing of the atmospheric general circulation.

875 To conclude, the application of the *High-Tune Explorer* on SCM/LES comparisons
876 allows us: (i) to quantify the parametric uncertainty at process level, (ii) to identify pa-
877 rameters which limit model performance, whatever their value, and should be replaced
878 by a more physical parameterization (i.e. when combining different cases, it may appear
879 that no value of a parameter is found acceptable for all cases and therefore suggests that
880 this parameter can not be kept constant but instead should depend on environmental
881 conditions), and (iii) to reduce the domain of acceptable values of free parameters used
882 in the final tuning of the global model.

883 We show indeed in Part II how the tool applied first to SCM/LES comparisons,
884 on a multicas e basis, can be used to reduce the range of acceptable values for the cal-
885 ibration of the complete 3D model configuration and considerably accelerate the resource
886 and time consumption for this step of model development. The final 3D tuning becomes
887 a part of the history matching process, by adding new metrics or constraints using the
888 exact same codes.

889 We believe that this tool is a breakthrough for model development as it allows us
890 to place the importance of the physical understanding of the processes at the heart of
891 model development, based on an extensive use of the SCM/LES comparison, whilst har-
892 nassing important techniques in machine learning and uncertainty quantification. We
893 advocate that the approach presented here leads to a well-defined strategy for calibra-
894 tion of the full model that may result in a significant acceleration in model improvement.

Table A1. List of the different LES runs of the Ayotte case used to determine the uncertainty around the reference

Name	Resolution	White noise	Turbulence	Diffusion
Name	Dx, Dz	Standard deviation (K)	length-scale	Timescale
Reference	50 m,nested <25 m	0.01 K	Deardorff length scale	1800 s
WhiteNoise	”	0.1 K	”	”
WhiteNoiseLL	”	0.5 K	”	”
Turb	”	”	size of the grid	”
Difshort	”	”	”	300 s
Diflong	”	”	”	7200 s
Dx	25 m, ”	”	”	”
Dz	”, nested <12.5 m	”	”	”

895 **Appendix A The different Large-Eddy Simulations**

896 Different simulations have been run with Meso-NH (Lac et al., 2018), varying the
897 resolution, domain size, turbulence formulation, intensity of the white noise introduced
898 at the first level and initial time to trigger turbulence, activation of subgrid condensa-
899 tion and changes in the microphysics scheme for the cloudy cases. The Table A1 lists
900 the different simulations of the Ayotte case used in section 4 to estimate the uncertainty
901 associated to the reference LES and the Table A2 lists the different simulations of the
902 ARMCU case used in section 3 to estimate the uncertainty associated to the reference
903 LES. The reference LES is highlighted in bold.

904 **Appendix B ARPEGE-Climat 6.3 and its turbulence parameteriza-** 905 **tion**

906 ARPEGE-Climat 6.3 is the atmospheric component of the CNRM-CM6-1 climate
907 model (Voltaire et al., 2019; Roehrig et al., 2020). It has 91 vertical levels, 15 of them
908 below 1500 m. The model time step is 15 minutes. Here, we use its SCM version and
909 focus on its representation of a clear convective boundary layer. To simulate the processes
910 involved in the boundary layer, the model combines a turbulence scheme with a mass-

Table A2. List of the different LES runs of the ARMCU case used to determine the uncertainty around the reference; the names indicated in the left column are those used in the legend of Figure 2

Name	Horizontal Resolution	Vertical Resolution	Domain side	Subgrid Condensation	Microphysics	Turbulence mixing length
12Dx25z25	25 m	25 m	12.8 km	No	Warm (ICE3)	Deardorff
6Dx25z25	”	”	6.4 km	”	”	”
6Dx40z25	40 m	25 m	6.4 km	”	”	”
6Dx40z40	40 m	40 m	6.4 km	”	”	”
6Dx25zvar	25 m	stretched grid	6.4 km	”	”	”
6Dx100z40	100 m	40 m	6.4 km	”	”	”
25Dx100z40	100 m	40 m	25.6 km	”	”	”
51Dx100z40	100 m	40 m	51.2 km	”	”	”
6DelDx25z25	25 m	25 m	6.4 km	”	”	$(Dx * Dy * Dz)^{1/3}$
6SbgDx25z25	25 m	25 m	6.4 km	Yes	”	Deardorff
6NprDx25z25	25 m	25 m	6.4 km	No	Only saturation adjustment	”

911 flux scheme, thus following the Eddy-Diffusivity Mass-Flux framework (e.g. Hourdin
 912 et al., 2002; Soares et al., 2004; Siebesma et al., 2007; Pergaud et al., 2009). The mass-
 913 flux scheme represents convection in a unified way from the clear convective boundary
 914 layer regime to the shallow cumulus and deep convection regimes (Piriou et al., 2007;
 915 Gueremy, 2011). In the illustration, we aim at analyzing the importance of the values
 916 of free parameters of the turbulence scheme on the simulation of an idealized clear bound-
 917 ary layer. A boundary-layer-top vertical entrainment is activated in the default version
 918 of ARPEGE-Climat 6.3 (see Roehrig et al., 2020)). For the sake of simplicity of the present
 919 illustration, and also because this parameterization is weakly active in the analyzed case,
 920 it is fully deactivated. Similar results are obtained when it is activated.

The turbulence scheme is based on Cuxart et al. (2000) which provides the verti-
 cal turbulent fluxes from which the turbulent source term is derived for the prognostic
 variables (see more details in Roehrig et al., 2020). The scheme relies on a prognostic
 equation of the grid-scale turbulence kinetic energy, \bar{e} :

$$\frac{\partial e}{\partial t} = \frac{-1}{\rho} \frac{\partial(\overline{\rho w' e'})}{\partial z} - (\overline{w' u'} \frac{\partial \bar{u}}{\partial z} + \overline{w' v'} \frac{\partial \bar{v}}{\partial z}) + \beta \overline{w' \theta'_{vl}} - \frac{\bar{e}^{3/2}}{L_\epsilon} \quad (\text{B1})$$

where the advection terms, the pressure fluctuations and the diffusion transport have
 been neglected. ρ is the air density, w the vertical velocity, u and v the zonal and merid-
 ional wind components, β is the buoyancy parameter (equal to $\frac{g}{\theta}$ with g the gravitational
 constant, θ being the potential temperature), θ_{vl} is the liquid virtual potential temper-
 ature and L_ϵ the dissipation length. Primes indicate fluctuations with respect to the grid-
 scale values indicated with overbars. The different turbulent vertical fluxes are diagnosed
 using \bar{e} following, for any variable φ :

$$\overline{w' \varphi'}(z) = -K_\varphi \frac{\partial \overline{\varphi}(z)}{\partial z} \quad (\text{B2})$$

with

$$K_\varphi = \sqrt{\bar{e}} L_m A_\varphi \Phi_\varphi \quad (\text{B3})$$

921 with Φ_φ a stability function also computed at each altitude (for more details see Cuxart
 922 et al. (2000)) and A_φ a free parameter. The mixing length, L_m , is computed following
 923 Bougeault and Lacarrere (1989); it consists in computing the vertical displacement an
 924 air parcel can travel upwards and downwards with its available turbulence kinetic en-
 925 ergy according to the thermal stratification. Also, L_ϵ in Eq. B1 is defined by $L_\epsilon = \frac{1}{A_\epsilon} \times$
 926 L_m with A_ϵ another free parameter.

927 **Acknowledgments**

928 This work received funding from grant HIGH-TUNE ANR-16-CE01-0010. It was sup-
 929 ported by the DEPHY2 project, funded by the French national program LEFE/INSU
 930 and the GDR-DEPHY. Daniel Williamson was funded by NERC grant: NE/N018486/1.
 931 Daniel Williamson and Victoria Volodina were funded by the Alan Turing Institute project
 932 “Uncertainty Quantification of multi-scale and multiphysics computer models: applica-
 933 tions to hazard and climate models” as part of the grant EP/N510129/1 made to the
 934 Alan Turing Institute by EPSRC. The authors would like to thank Ignacio Lopez-Gomez
 935 and an anonymous reviewer for their constructive comments. We also thank S Barbier,
 936 T Costablos, S Nicolau and S Richet for their work during a short internship on that sub-
 937 ject. Beyond the presentation of a new approach that we think could constitute a break
 938 through in climate model improvement, we intend to provide a tool for the climate com-
 939 munity. All the programs, scripts and reference LES are publicly available via a Subver-
 940 sion through ”svn checkout <https://svn.lmd.jussieu.fr/HighTune>”. Note, however, that
 941 this tool is a new research tool, and, as such, is still evolving. The code, the SCM runs
 942 and the LES used to produce Experiment 1 is available at doi: XXXX - will be provided
 943 during proofs.

944 **References**

- 945 Ahmat Younous, A.-L., Roehrig, R., Beau, I., & Douville, H. (2018). Single-column
 946 modeling of convection during the cindy2011/dynamo field campaign with
 947 the cnrm climate model version 6. *Journal of Advances in Modeling Earth*
 948 *Systems*, 10, 578–602.
- 949 Andrianakis, I., McCreesh, N., Vernon, I., McKinley, T. J., Oakley, J. E., Nsub-
 950 uga, R. N., ... White, R. G. (2017). Efficient history matching of a high
 951 dimensional individual-based hiv transmission model. *SIAM/ASA Journal on*
 952 *Uncertainty Quantification*, 5(1), 694–719.
- 953 Ayotte, K. W., Sullivan, P. P., Andren, A., Doney, S. C., Holtslag, A. A., Large,
 954 W. G., ... Wyngaard, J. C. (1996). An evaluation of neutral and convective
 955 planetary boundary-layer parameterizations relative to large eddy simulations.
 956 *Boundary-layer Meteorol.*, 79, 131–175.
- 957 Bastidas, L. A., Hogue, T. S., Sorooshian, S., Gupta, H. V., & J, S. W. (2006).
 958 Parameter sensitivity analysis for different complexity land surface models

- 959 using multicriteria methods. *Journal of Geophysical Research*, *111*. doi:
 960 10.1029/2005JD006377
- 961 Bellprat, ., Kotlarski, S., Luthi, D., & Schar, C. (2012). Objective calibration of re-
 962 gional climate models. *Journal of Geophysical Research*, *117*. doi: 10.1029/
 963 2012JD018262
- 964 Bony, S., Stevens, B., Frierson, D. M. W., Jakob, C., Kageyama, M., Pincus, R.,
 965 ... Webb, M. J. (2015, April). Clouds, cirulation and climate sensitiv-
 966 ity. *Nature Geoscience*, *8*(20), L20806. (WOS:000233104900005) doi:
 967 10.1038/NGEO2398
- 968 Bougeault, P., & Lacarrere, P. (1989). Parameterization of orography induced turbu-
 969 lence in a mesobeta-scale model. *Mon. Wea. Rev.*, *117*, 1872–1890.
- 970 Bouniol, D., Roca, R., Fiolleau, T., & Poan, E. (2016). Macrophysical, microphysical
 971 and radiative properties of tropical mesoscale convective systems over their life
 972 cycle. *Journal of Climate*, *29*. doi: 10.1175/JCLI-D-15-0551.1
- 973 Brenowitz, N. D., & Bretherton, C. S. (2018, June). Prognostic validation of a
 974 neural network unified physics parameterization. *Geophysical Research Letters*,
 975 *45*(12), 6289–6298. (WOS:000438499100052) doi: 10.1029/2018GL078510
- 976 Brient, F., Couvreur, F., Villefranque, N., Rio, C., & Honnert, R. (2019). Object-
 977 oriented identification of coherent structures in large eddy simulations: Im-
 978 portance of downdrafts in stratocumulus. *Geophysical Research Letters*, *46*,
 979 2854–2864.
- 980 Brown, A. R. (1999, January). The sensitivity of large-eddy simulations of shallow
 981 cumulus convection to resolution and subgrid model. *Quarterly Journal of the*
 982 *Royal Meteorological Society*, *125*(554), 469–482. (WOS:000079350100004) doi:
 983 10.1002/qj.49712555405
- 984 Brown, A. R., Cederwall, R. T., Chlond, A., Duynkerke, P. G., Golaz, M., J.
 985 C. Khairoutdinov, Lewellen, D. C., ... Stevens, B. (2002). Large-eddy simu-
 986 lation of the diurnal cycle of shallow cumulus convection over land. *Q. J. R.*
 987 *Meteorol. Soc.*, *128*, 1075–1093.
- 988 Browning, K., Betts, A., Jonas, P., Kershaw, R., Manton, M., Mason, P.,
 989 ... Simpson, J. (1993, March). The GEWEX Cloud System Study
 990 (GCSS). *Bulletin of the American Meteorological Society*, *74*(3), 387–399.
 991 (WOS:A1993KU53500004)

- 992 Brynjarsdóttir, J., & O’Hagan, A. (2014). Learning about physical parameters: The
 993 importance of model discrepancy. *Inverse Problems*, *30*(11), 114007.
- 994 Caldwell, P., & Bretherton, C. S. (2009, January). Response of a Subtropi-
 995 cal Stratocumulus-Capped Mixed Layer to Climate and Aerosol Changes.
 996 *Journal of Climate*, *22*(1), 20–38. (WOS:000262329100002) doi: 10.1175/
 997 2008JCLI1967.1
- 998 Carpenter, B., & Coauthors. (2017, May). Stan: A probabilistic programming lan-
 999 guage. *Journal of Statistical Software*, *76*(1), 00–00. doi: 10.18637/jss.v076
 1000 .i01
- 1001 Chinita, M. J., Matheou, G., & Teixeira, J. (2018). A joint probability density
 1002 based decomposition of turbulence in the atmospheric boundary layer. *Monthly*
 1003 *Weather Review*, *146*, 503–523.
- 1004 Couvreux, F., Guichard, F., Redelsperger, J. L., Kiemle, C., Masson, V., Lafore,
 1005 J. P., & Flamant, C. (2005, October). Water-vapour variability within a
 1006 convective boundary-layer assessed by large-eddy simulations and IHOP_2002
 1007 observations. *Quarterly Journal of the Royal Meteorological Society*, *131*(611),
 1008 2665–2693. (WOS:000233475900005) doi: 10.1256/qj.04.167
- 1009 Couvreux, F., Hourdin, F., & Rio, C. (2010, March). Resolved Versus Parametrized
 1010 Boundary-Layer Plumes. Part I: A Parametrization-Oriented Conditional
 1011 Sampling in Large-Eddy Simulations. *Boundary-Layer Meteorology*, *134*(3),
 1012 441–458. (WOS:000274013600004) doi: 10.1007/s10546-009-9456-5
- 1013 Craig, P. S., Goldstein, M., Seheult, A., & Smith, J. (1996). Bayes linear strategies
 1014 for matching hydrocarbon reservoir history. *Bayesian statistics*, *5*, 69–95.
- 1015 Cuxart, J., Bougeault, P., & Redelsperger, J.-L. (2000). A turbulence scheme allow-
 1016 ing for mesoscale and large-eddy simulations. *Q. J. R. Meteorol. Soc.*, *126*, 1–
 1017 30.
- 1018 de Roode, S. R., Sandu, I., Van Der Dussen, J. J., Ackerman, A. S., Blossey, P.,
 1019 Jarecka, D., . . . Stevens, B. (2016). Large-eddy simulations of euclipse-gass
 1020 lagrangian stratocumulus-to-cumulus transitions: Mean state, turbulence, and
 1021 decoupling. *Journal of the Atmospheric Sciences*, *73*, 2485–2508.
- 1022 Duan, Q., Di, Z., Quan, J., Wang, C., Gong, W., Gan, Y., . . . S, F. (2017). Au-
 1023 tomatic model calibration: A new way to improve numerical weather fore-
 1024 casting. *Bulletin of American Meteorological Society*, *98*, 959–970. doi:

1025 10.1175/BAMS-D-15-00104.1

1026 Duynkerke, P. G., de Roode, S. R., van Zanten, M. C., Calvo, J., Cuxart, J., &
 1027 Cheinet, S. (2004). Observations and numerical simulations of the diurnal
 1028 cycle of the eurocs stratocumulus case. *Quarterly Journal of the Royal Meteorological Society*, *130*, 3269–3296.
 1029

1030 Flato, J., G. and Marotzke, Abiodun, B., Braconnot, P., Chou, S., Collins, W., Cox,
 1031 P., . . . Rummukaine, M. (2013, August). Evaluation of climate models. *Climate Change 2013: The Physical Science Basis. Contribution of Working*
 1032 *Group I to the Fifth Assessment Report of the Intergovernmental Panel on*
 1033 *Climate Change*.
 1034

1035 Gelman, A. (2006). Prior distributions for variance parameters in hierarchical mod-
 1036 els. *Bayesian Analysis*, *1*, 515–534.

1037 Gentine, P., Pritchard, M., Rasp, S., Reinaudi, G., & Yacalis, G. (2018, May). Could
 1038 machine learning break the convection parameterization deadlock? *Geophysical*
 1039 *Research Letters*, *45*, 5742–5751. doi: 10.1029/2018GL078202

1040 Gettelman, A., Truesdale, J., Bacmeister, J., Caldwell, P., Neale, R., & Bogenschutz,
 1041 P. (2019, May). The single column atmosphere model version 6 (SCAM6): Not
 1042 a scam but a tool for model evaluation and development. *Journal of Advances*
 1043 *in modeling earth systems*, *11*, 1381–1401. doi: 10.1029/2018MS001578

1044 Golaz, J.-C., Horowitz, L. W., & Levy, H. (2013, May). Cloud tuning in a cou-
 1045 pled climate model: Impact on 20th century warming. *Geophysical Research*
 1046 *Letters*, *40*(10), 2246–2251. (WOS:000328840200064) doi: 10.1002/grl.50232

1047 Golaz, J. C., Larson, V. E., & Cotton, W. R. (2002, December). A PDF-based
 1048 model for boundary layer clouds. Part II: Model results. *Journal of the Atmo-*
 1049 *spheric Sciences*, *59*(24), 3552–3571. (WOS:000179629800007) doi: 10.1175/
 1050 1520-0469(2002)059(3552:APBMFB)2.0.CO;2

1051 Grabowski, W. W. (2016, June). Towards global large-eddy simulation: super pa-
 1052 rameterization revisited. *Journal of the Meteorological Society of Japan*, *94*(4),
 1053 L20806. doi: 10.2151/jmsj.2016-017

1054 Gueremy, J. F. (2011, August). A continuous buoyancy based convection scheme:
 1055 one- and three-dimensional validation. *Tellus Series a-Dynamic Meteorology*
 1056 *and Oceanography*, *63*(4), 687–706. (WOS:000292864500004) doi: 10.1111/j
 1057 .1600-0870.2011.00521.x

- 1058 Guichard, F., & Couvreux, F. (2017). A short review of numerical cloud-
 1059 resolving models. *Tellus Dyn. Meteorol. Oceanogr.*, *69*, 1945–1960. doi:
 1060 10.1080/16000870.2017.1373578
- 1061 Guo, Z., Wong, T. S., Larson, V. E., Ghan, S., Ovchinnikov, M., Bogenschutz,
 1062 P. A., ... Zhou, T. (2014). A sensitivity analysis of cloud properties
 1063 to clubb parameters in the single-column community atmosphere model
 1064 (scam5). *Journal of Advances in Modeling Earth Systems*, *6*, 829–858. doi:
 1065 10.1002/2014MS000315
- 1066 Heus, T., & Jonker, H. J. J. (2008, March). Subsiding shells around shallow
 1067 cumulus clouds. *Journal of the Atmospheric Sciences*, *65*(3), 1003–1018.
 1068 (WOS:000254356600016) doi: 10.1175/2007JAS2322.1
- 1069 Heus, T., Pols, C. F. J., Jonker, H. J. J., Van den Akker, H. E. A., & Lenschow,
 1070 D. H. (2009, January). Observational validation of the compensating mass
 1071 flux through the shell around cumulus clouds. *Quarterly Journal of the Royal*
 1072 *Meteorological Society*, *135*(638), 101–112. (WOS:000265374900008) doi:
 1073 10.1002/qj.358
- 1074 Holtslag, A. A. M., Svensson, G., Baas, P., Basu, S., Beare, B., Beljaars, A. C. M.,
 1075 ... Van de Wiel, B. J. H. (2013, November). STABLE ATMOSPHERIC
 1076 BOUNDARY LAYERS AND DIURNAL CYCLES Challenges for Weather
 1077 and Climate Models. *Bulletin of the American Meteorological Society*, *94*(11),
 1078 1691–1706. (WOS:000327926700007) doi: 10.1175/BAMS-D-11-00187.1
- 1079 Hourdin, F., Couvreux, F., & Menut, L. (2002, March). Parameterization
 1080 of the dry convective boundary layer based on a mass flux representa-
 1081 tion of thermals. *Journal of the Atmospheric Sciences*, *59*(6), 1105–
 1082 1123. (WOS:000174019900006) doi: 10.1175/1520-0469(2002)059<1105:
 1083 POTDCB>2.0.CO;2
- 1084 Hourdin, F., Grandpeix, J.-Y., Rio, C., Bony, S., Jam, A., Cheruy, F., ... Roehrig,
 1085 R. (2013, May). LMDZ5b: the atmospheric component of the IPSL cli-
 1086 mate model with revisited parameterizations for clouds and convection.
 1087 *Climate Dynamics*, *40*(9-10), 2193–2222. (WOS:000318278700005) doi:
 1088 10.1007/s00382-012-1343-y
- 1089 Hourdin, F., Mauritsen, T., Gettelman, A., Golaz, J.-C., Balaji, V., Duan, Q., ...
 1090 Williamson, D. (2017, March). The Art and Science of Climate Model Tuning.

- 1091 *Bull. Am. Meteorol. Soc.*, *98*, 589-602. doi: 10.1175/BAMS-D-15-00135.1
- 1092 Hourdin, F., Rio, C., Grandpeix, J.-Y., Madeleine, J.-B., Cheruy, F., Rochetin,
 1093 N., . . . Ghattas, J. (2020, June). LMDZ6A: the atmospheric component
 1094 of the ipsl climate model with improved and better tuned physics. *Journal*
 1095 *of Advances in modeling earth systems, accepted for publication, ??-??* doi:
 1096 10.1029/2019MS001892
- 1097 Jakob, C. (2010, July). ACCELERATING PROGRESS IN GLOBAL ATMO-
 1098 SPHERIC MODEL DEVELOPMENT THROUGH IMPROVED PARAM-
 1099 ETERIZATIONS Challenges, Opportunities, and Strategies. *Bulletin of the*
 1100 *American Meteorological Society*, *91*(7), 869–+. (WOS:000280758700003) doi:
 1101 10.1175/2009BAMS2898.1
- 1102 Jam, A., Hourdin, F., Rio, C., & Couvreux, F. (2013, June). Resolved Versus
 1103 Parametrized Boundary-Layer Plumes. Part III: Derivation of a Statistical
 1104 Scheme for Cumulus Clouds. *Boundary-Layer Meteorology*, *147*(3), 421–441.
 1105 (WOS:000319475000004) doi: 10.1007/s10546-012-9789-3
- 1106 Jiang, J. H., Su, H., Zhai, C., Perun, V., Del Genio, A., Nazarenko, L. S., . . . L,
 1107 S. G. (2012, July). Evaluation of cloud and water vapor simulations in CMIP5
 1108 climate models using nasa “a-train” satellite observations. *Journal of Geophysic-*
 1109 *al Research*, *117*, 1–24. doi: 10.1029/2011JD017237
- 1110 Johnson, J. S., Cui, Z., Lee, L. A., Gosling, J. P., Blyth, A. M., & Carslaw, K. S.
 1111 (2015). Evaluating uncertainty in convective cloud microphysics using statisti-
 1112 cal emulation. *Journal of Advances in Modeling Earth Systems*, *7*, 162–187.
- 1113 Kennedy, M. C., & O’Hagan, A. (2001, aug). Bayesian calibration of computer mod-
 1114 els. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*,
 1115 *63*(3), 425–464. Retrieved from [http://doi.wiley.com/10.1111/1467-9868](http://doi.wiley.com/10.1111/1467-9868.00294)
 1116 [.00294](http://doi.wiley.com/10.1111/1467-9868.00294) doi: 10.1111/1467-9868.00294
- 1117 Khairoutdinov, M., Randall, D., & DeMott, C. (2005, July). Simulations of the
 1118 atmospheric general circulation using a cloud-resolving model as a superpa-
 1119 rameterization of physical processes. *Journal of the Atmospheric Sciences*,
 1120 *62*(7), 2136–2154. (WOS:000230962800006) doi: 10.1175/JAS3453.1
- 1121 Klein, S. A., Hall, A., R., N. J., & Robert, P. (2017, October). Low-cloud feedbacks
 1122 from cloud-controlling factors: a review. *Survey of Geophysics*, *38*(10), 1307–
 1123 1329. doi: 10.1007/s10712-017-9433-3

- 1124 Krasnopolsky, V. M., Fox-Rabinovitz, M. S., & Belochitski, A. A. (2013, March).
 1125 Using ensemble of neural networks to learn stochastic convection parameteriza-
 1126 tions for climate and numerical weather prediction models from data simulated
 1127 by a cloud resolving model. *Advances in Artificial Neural Systems*, 203(3), 13.
 1128 doi: 10.1155/2013/485913
- 1129 Kumar, V. V., Jakob, C., Protat, A., Williams, C. R., & May, P. T. (2015, May).
 1130 Mass-Flux Characteristics of Tropical Cumulus Clouds from Wind Profiler Ob-
 1131 servations at Darwin, Australia. *Journal of the Atmospheric Sciences*, 72(5),
 1132 1837–1855. (WOS:000353840100009) doi: 10.1175/JAS-D-14-0259.1
- 1133 Lac, C., Chaboureaud, J.-P., Masson, V., Pinty, J.-P., Tulet, P., Escobar, J., ...
 1134 Wautelet, P. (2018, January). Overview of the Meso-NH model version 5.4 and
 1135 its applications. *Geosci. Model Dev. Discuss.*, 2018, 1–66. Retrieved 2018-03-
 1136 26, from <https://www.geosci-model-dev-discuss.net/gmd-2017-297/> doi:
 1137 10.5194/gmd-2017-297
- 1138 Loeppky, J. L., Sacks, J., & Welch, W. J. (2009). Choosing the sample size of a com-
 1139 puter experiment: A practical guide. *Technometrics*, 51, 366–376. doi: 10
 1140 .1198/TECH.2009.08040
- 1141 Masunaga, H. (2012, November). Short-Term versus Climatological Relationship
 1142 between Precipitation and Tropospheric Humidity. *Journal of Climate*, 25(22),
 1143 7983–7990. (WOS:000311034300013) doi: 10.1175/JCLI-D-12-00037.1
- 1144 Masunaga, H., & Luo, Z. L. (2016). Convective and large-scale mass flux profiles
 1145 over tropical oceans determined from synergetic analysis of a suite of satel-
 1146 lite observations. *Journal of Geophysical Research*, 121, 7958–7974. doi:
 1147 10.1002/2016JD024753
- 1148 Matheou, G., Chung, D., Nuijens, L., Stevens, B., & Teixeira, J. (2011, Septem-
 1149 ber). On the Fidelity of Large-Eddy Simulation of Shallow Precipitat-
 1150 ing Cumulus Convection. *Monthly Weather Review*, 139(9), 2918–2939.
 1151 (WOS:000294932100014) doi: 10.1175/2011MWR3599.1
- 1152 Mauritsen, T., Stevens, B., Roeckner, E., Crueger, T., Esch, M., Giorgetta, M., ...
 1153 Tomassini, L. (2012, August). Tuning the climate of a global model. *Journal*
 1154 *of Advances in Modeling Earth Systems*, 4, M00A01. (WOS:000307467200001)
 1155 doi: 10.1029/2012MS000154
- 1156 McNeall, D., Williams, J., Betts, R., Booth, B., Challenor, P., Good, P., & Wilt-

- 1157 shire, A. (2019). Correcting a bias in a climate model with an augmented
 1158 emulator. *Geoscientific Model Development Discussions, 2019*, 1–37. Retrieved
 1159 from <https://www.geosci-model-dev-discuss.net/gmd-2019-171/> doi:
 1160 10.5194/gmd-2019-171
- 1161 Nam, C., Bony, S., Dufresne, J.-L., & Chepfer, H. (2012, November). The 'too few,
 1162 too bright' tropical low-cloud problem in CMIP5 models. *Geophysical Research*
 1163 *Letters, 39*, L21801. (WOS:000310690600003) doi: 10.1029/2012GL053421
- 1164 Neggers, R. A. J. (2009, June). A Dual Mass Flux Framework for Boundary Layer
 1165 Convection. Part II: Clouds. *Journal of the Atmospheric Sciences, 66*(6),
 1166 1489–1506. (WOS:000267263300002) doi: 10.1175/2008JAS2636.1
- 1167 Neggers, R. A. J. (2015, December). Attributing the behavior of low-level clouds in
 1168 large-scale models to subgrid-scale parameterizations. *Journal of Advances in*
 1169 *modeling earth systems, 7*(4), 2029–2043. doi: 10.1002/2015MS000503
- 1170 Neggers, R. A. J., Ackerman, A. S., Angevine, W. M., Bazile, E., Beau, I., Blossey,
 1171 P. N., ... Heus, T. (2017, October). Single-column model simulations of
 1172 subtropical marine boundary-layer cloud transitions under weakening inver-
 1173 sions. *Journal of Advances in modeling earth systems, 9*(6), 2385–2412. doi:
 1174 10.1002/2017MS001064
- 1175 Neggers, R. A. J., Duynkerke, P. G., & Rodts, S. M. A. (2003, July). Shallow cu-
 1176 mulus convection: A validation of large-eddy simulation against aircraft and
 1177 Landsat observations. *Quarterly Journal of the Royal Meteorological Society,*
 1178 *129*(593), 2671–2696. (WOS:000185187600011) doi: 10.1256/qj.02.93
- 1179 Neggers, R. A. J., Jonker, H. J., & Siebesma, P. (2003). Statistics of cumulus cloud
 1180 populations in large-eddy simulations. *J. Atmos. Sci., 60*, 1060–1074.
- 1181 Neggers, R. A. J., Siebesma, A. P., & Heus, T. (2012, September). Continu-
 1182 ous single-column model evaluation at a permanent meteorological super-
 1183 site. *Bulletin of the American Meteorological Society, 93*(9), 1389–1400.
 1184 (WOS:000309056400006) doi: 10.1175/BAMS-D-11-00162.1
- 1185 Neggers, R. A. J., Siebesma, A. P., Lenderink, G., & Holtslag, A. A. M. (2004,
 1186 November). An evaluation of mass flux closures for diurnal cycles
 1187 of shallow cumulus. *Monthly Weather Review, 132*(11), 2525–2538.
 1188 (WOS:000225098900001) doi: 10.1175/MWR2776.1
- 1189 Neggers, R. A. J., Siebesma, P., & J, J. H. J. (2002). A multiparcel model for shal-

- 1190 low cumulus convection. *J. Atmos. Sci.*, *59*, 1655–1668.
- 1191 Nuijens, L., Medeiros, B., Sandu, I., & Ahlgrimm, M. (2015, December). Observed
1192 and modeled patterns of covariability between low-level cloudiness and the
1193 structure of the trade-wind layer. *Journal of Advances in Modeling Earth Sys-*
1194 *tems*, *7*(4), 1741–1764. (WOS:000368739800013) doi: 10.1002/2015MS000483
- 1195 Oakley, J. E., & O’Hagan, A. (2004, December). Probabilistic sensitivity analysis
1196 of complex models: a bayesian approach. *Royal Statistical Society*, *66*(3), 751–
1197 769.
- 1198 Parishani, H., Pritchard, M., Bretherton, C., Wyant, M., & Khairoutdinov, M.
1199 (2017, May). Toward low-cloud-permitting cloud superparameterization with
1200 explicit boundary layer turbulence. *Journal of Advances in modeling earth*
1201 *systems*, *9*, 1542–1571. doi: 10.1002/2018MS001409
- 1202 Pergaud, J., Masson, V., Malardel, S., & Couvreux, F. (2009, July). A Param-
1203 eterization of Dry Thermals and Shallow Cumuli for Mesoscale Numeri-
1204 cal Weather Prediction. *Boundary-Layer Meteorology*, *132*(1), 83–106.
1205 (WOS:000267029600006) doi: 10.1007/s10546-009-9388-0
- 1206 Piriou, J.-M., Redelsperger, J.-L., Geleyn, J.-F., Lafore, J.-P., & Guichard, F. (2007,
1207 November). An approach for convective parameterization with memory: Sep-
1208 arating microphysics and transport in grid-scale equations. *Journal of the At-*
1209 *mospheric Sciences*, *64*(11), 4127–4139. (WOS:000251283000025) doi: 10
1210 .1175/2007JAS2144.1
- 1211 Pressel, K. G., Mishra, S., Schneider, T., Kaul, C. M., & Tan, Z. (2017, June). Nu-
1212 merics and subgrid-scale modeling in large eddy simulations of stratocumulus
1213 clouds. *Journal of Advances in Modeling Earth Systems*, *9*(2), 1342–1365. doi:
1214 10.1002/2016MS000778
- 1215 Pukelsheim, F. (1994, February). The three sigma rule. *The American Statistician*,
1216 *48*, 88–91.
- 1217 Randall, D., Khairoutdinov, M., Arakawa, A., & Grabowski, W. (2003, November).
1218 Breaking the cloud parameterization deadlock. *Bulletin of the American Mete-*
1219 *orological Society*, *84*(11), 1547–1564. (WOS:000187163900019) doi: 10.1175/
1220 BAMS-84-11-1547
- 1221 Randall, D., Xu, K., Somerville, R., & Iacobellis, S. (1996, August). Single-column
1222 models and cloud ensemble models as links between observations and climate

- 1223 models. *Journal of Climate*, 9(8), 1683–1697. (WOS:A1996VG92100002) doi:
 1224 10.1175/1520-0442(1996)009<1683:SCMACE>2.0.CO;2
- 1225 Richter, I. (2015, February). Climate model biases in the eastern tropical oceans:
 1226 Causes, impacts and ways forward. *Wiley Interdisciplinary Reviews:Climate*
 1227 *Change*, 6(3), 345–358.
- 1228 Rio, C., Del Genio, A. D., & F, H. (2019). Ongoing breakthroughs in convective pa-
 1229 rameterization. *Current Climate Change Reports*, 5, 95–111.
- 1230 Rio, C., & Hourdin, F. (2008, February). A thermal plume model for the convective
 1231 boundary layer: Representation of cumulus clouds. *Journal of the Atmospheric*
 1232 *Sciences*, 65(2), 407–425. (WOS:000253406600007) doi: 10.1175/2007JAS2256
 1233 .1
- 1234 Rio, C., Hourdin, F., Couvreux, F., & Jam, A. (2010, June). Resolved Versus
 1235 Parametrized Boundary-Layer Plumes. Part II: Continuous Formulations of
 1236 Mixing Rates for Mass-Flux Schemes. *Boundary-Layer Meteorology*, 135(3),
 1237 469–483. (WOS:000277635800007) doi: 10.1007/s10546-010-9478-z
- 1238 Rochetin, N., Couvreux, F., Grandpeix, J.-Y., & Rio, C. (2014, February). Deep
 1239 Convection Triggering by Boundary Layer Thermals. Part I: LES Analysis
 1240 and Stochastic Triggering Formulation. *Journal of the Atmospheric Sciences*,
 1241 71(2), 496–514. (WOS:000335491400003) doi: 10.1175/JAS-D-12-0336.1
- 1242 Roehrig, R., Beau, I., Saint-Martin, D., Alias, A., Colin, J., Decharme, B., . . .
 1243 S en esi, S. (2020, June). The CNRM global atmosphere model ARPEGE-
 1244 Climat 6.3: description and evaluation. *Journal of Advances in modeling earth*
 1245 *systems, accepted for publication, ??-??* doi: 10.1029/2020MS002075
- 1246 Rougier, J., Sexton, D. M. H., Murphy, J. M., & Stainforth, D. (2009). Analyz-
 1247 ing the climate sensitivity of the hadsm3 climate model using ensembles from
 1248 different but related experiments. *Journal of Climate*, 22, 3540–3557. doi:
 1249 10.1175/2008JCLI2533.1
- 1250 Saltelli, A. J. (2002, December). Making best use of model evaluations to compute
 1251 sensitivity indices. *Computer Physics Communications*, 145(2), 280–297.
- 1252 Salter, J. M., Williamson, D. B., Scinocca, J., & Kharin, V. (2019, October).
 1253 Uncertainty quantification for computer models with spatial output using
 1254 calibration-optimal bases. *Journal of the American statistical association*,
 1255 114(528), 1800–1814. doi: 10.1080/01621459.2018.1514306

- 1256 Sandu, I., Beljaars, A., Bechtold, P., Mauritsen, T., & Balsamo, G. (2013). Why
 1257 is it so difficult to represent stably stratified conditions in numerical weather
 1258 prediction (NWP) models? *Journal of Advances in Modeling Earth Systems*,
 1259 5(2), 117–133. Retrieved from <http://dx.doi.org/10.1002/jame.20013> doi:
 1260 10.1002/jame.20013
- 1261 Satoh, M., Matsuno, T., Tomita, H., Miura, H., Nasuno, T., & Iga, S. (2008,
 1262 March). Nonhydrostatic icosahedral atmospheric model (NICAM) for global
 1263 cloud resolving simulations. *Journal of Computational Physics*, 227(7), 3486–
 1264 3514. (WOS:000255005900005) doi: 10.1016/j.jcp.2007.02.006
- 1265 Satoh, M., Stevens, B., Judt, F., Khairoutdinov, M., Lin, S.-J., Putman, W., & P.,
 1266 D. (2019, September). Global cloud-resolving models. *Current Climate Change*
 1267 *Reports*, 5(3), 172–184. doi: 10.1007/s40641-019-00131-0
- 1268 Schmidt, G. A., Bader, D., Donner, L. J., Elsaesser, G. S., Golaz, J.-C., Hannay, C.,
 1269 ... Saha, S. (2017, September). Practice and philosophy of climate model tun-
 1270 ing across six US modeling centers. *Geoscientific Model Development*, 10(9),
 1271 3207–3223. (WOS:000409053900001) doi: 10.5194/gmd-10-3207-2017
- 1272 Schneider, T., Lan, T., Stuart, A., & Teixeira, J. (2017, December). Earth sys-
 1273 tem modeling 2.0: A blueprint for models that learn from observations and
 1274 targeted high-resolution simulations. *Geophysical Research Letters*, 44, 12396–
 1275 12417. doi: 10.1002/2017GL076101
- 1276 Sexton, D. M., Murphy, J. M., Collins, M., & Webb, M. J. (2011, October). Multi-
 1277 variate probabilistic projections using imperfect climate models part i: outline
 1278 of methodology. *Climate Dynamics*, 38(1), 2513–2542.
- 1279 Siebesma, A. P., Bretherton, C. S., Brown, A., Chlond, A., Cuxart, J., Duynkerke,
 1280 P. G., ... Stevens, D. E. (2003). A large eddy simulation intercomparison
 1281 study of shallow cumulus convection. *J. Atmos. Sci.*, 60, 1201–1219.
- 1282 Siebesma, A. P., & Cuijpers, J. W. M. (1995). Evaluation of parametric assumptions
 1283 for shallow cumulus convection. *J. Atmos. Sci.*, 52, 650–666.
- 1284 Siebesma, A. P., Soares, P. M. M., & Teixeira, J. (2007, April). A combined eddy-
 1285 diffusivity mass-flux approach for the convective boundary layer. *Journal of*
 1286 *the Atmospheric Sciences*, 64(4), 1230–1248. (WOS:000245742600011) doi:
 1287 10.1175/JAS3888.1
- 1288 Soares, P. M. M., Miranda, P. M. A., Siebesma, A. P., & Teixeira, J. (2004). An

- 1289 eddy-diffusivity/mass-flux parameterization for dry and shallow cumulus con-
 1290 vection. *Q. J. R. Meteorol. Soc.*, *130*(604), 3365–3383.
- 1291 Stevens, B., Moeng, C. H., Ackerman, A. S., Bretherton, C. S., Chlond, A.,
 1292 De Roode, S., . . . Zhu, P. (2005, June). Evaluation of large-Eddy simulations
 1293 via observations of nocturnal marine stratocumulus. *Monthly Weather Review*,
 1294 *133*(6), 1443–1462. (WOS:000230028000003) doi: 10.1175/MWR2930.1
- 1295 Stevens, B., Satoh, M., Auger, L., Bierchamp, J., Bretherton, C. S., Chen, X., . . .
 1296 Chou, L. (2019). Dyamond: the dynamics of the atmospheric general circu-
 1297 lation modeled on non-hydrostatic domains. *Progress in Earth and Planetary*
 1298 *Science*, *6*, 1–17. doi: 10.1186/s40645-019-0304-z
- 1299 Sullivan, P. P., & Patton, E. G. (2011, October). The Effect of Mesh Resolution
 1300 on Convective Boundary Layer Statistics and Structures Generated by Large-
 1301 Eddy Simulation. *Journal of the Atmospheric Sciences*, *68*(10), 2395–2415.
 1302 (WOS:000296034700014) doi: 10.1175/JAS-D-10-05010.1
- 1303 Suselj, K., Kurowski, M. J., & Teixeira, J. (2019, August). A unified eddy-
 1304 diffusivity/mass-flux approach for modeling atmospheric convection. *Journal of*
 1305 *the Atmospheric Sciences*, 2505–2537.
- 1306 Suselj, K., Teixeira, J., & Chung, D. (2013, July). A Unified Model for Moist Con-
 1307 vective Boundary Layers Based on a Stochastic Eddy-Diffusivity/Mass-Flux
 1308 Parameterization. *Journal of the Atmospheric Sciences*, *70*(7), 1929–1953.
 1309 (WOS:000322125600005) doi: 10.1175/JAS-D-12-0106.1
- 1310 Tan, Z., Kaul, C. M., Pressel, G., K, Cohen, Y., Schneider, T., & Teixeira, J. (2018,
 1311 March). An extended eddy-diffusivity mass-flux scheme for unified representa-
 1312 tion of subgrid scale turbulence and convection. *Journal of Advances in Model-*
 1313 *ing Earth Systems*, 770–800.
- 1314 vanZanten, M. C., Stevens, B., Nuijens, L., Siebesma, A. P., Ackerman, A. S., Bur-
 1315 net, F., . . . Wyszogrodzki, A. (2011). Controls on precipitation and cloudiness
 1316 in simulations of trade-wind cumulus as observed during RICO. *Journal of*
 1317 *Advances in Modeling Earth Systems*, *3*, M06001. (WOS:000303198400003)
 1318 doi: 10.1029/2011MS000056
- 1319 Vernon, I., Goldstein, M., & Bower, R. (2010). Galaxy formation: a bayesian uncer-
 1320 tainty analysis. *Bayesian Analytics*, *5*, 619–846.
- 1321 Villefranque, N., Fournier, R., Couvreur, F., Blanco, S., Eymet, V., Forest, V., &

- 1322 Tregan, J. M. (2019). A path-tracing monte carlo library for 3-d radiative
 1323 transfer in highly resolved cloudy atmospheres. *Journal of Advances in Model-*
 1324 *ing Earth Systems*, *11*, 2449–2473.
- 1325 Voldoire, A., Saint-Martin, D., Senesi, S., Decharme, B., Alias, A., Chevallier, M.,
 1326 ... Waldman, W., R (2019, August). Evaluation of CMIP6 deck experiments
 1327 with CNRM-CM6-1. *Journal of Advances in Modeling Earth Systems*, *11*(7),
 1328 2177–2213. doi: 10.1029/2019MS001683
- 1329 Volodina, V. (2020). *Uncertainty quantification for complex computer models with*
 1330 *nonstationary output. bayesian optimal design for iterative refocussing* (Unpub-
 1331 lished doctoral dissertation). University of Exeter.
- 1332 Volodina, V., & Williamson, D. (2020, January). Diagnostics-driven nonstationary
 1333 emulators using kernel mixtures. *Journal of Uncertainty Quantification*, *8*(1),
 1334 1–26.
- 1335 Wang, H., & Feingold, G. (2009, November). Modeling Mesoscale Cellular Struc-
 1336 tures and Drizzle in Marine Stratocumulus. Part I: Impact of Drizzle on the
 1337 Formation and Evolution of Open Cells. *Journal of the Atmospheric Sciences*,
 1338 *66*(11), 3237–3256. (WOS:000271689700001) doi: 10.1175/2009JAS3022.1
- 1339 Williamson, D. (2015, June). Exploratory ensemble designs for environmental
 1340 models using k-extended Latin Hypercubes. *Environmetrics*, *26*(4), 268–283.
 1341 (WOS:000353380200003) doi: 10.1002/env.2335
- 1342 Williamson, D., Blaker, A. T., Hampton, C., & Salter, J. (2015, September). Iden-
 1343 tifying and removing structural biases in climate models with history match-
 1344 ing. *Climate Dynamics*, *45*(5-6), 1299–1324. (WOS:000360507700010) doi:
 1345 10.1007/s00382-014-2378-z
- 1346 Williamson, D., Blaker, A. T., & Sinha, B. (2017, April). Tuning without over-
 1347 tuning: parametric uncertainty quantification for the NEMO ocean model.
 1348 *Geoscientific Model Development*, *10*(4), 1789–1816. (WOS:000400181200002)
 1349 doi: 10.5194/gmd-10-1789-2017
- 1350 Williamson, D., Goldstein, M., Allison, L., Blaker, A., Challenor, P., Jackson, L.,
 1351 & Yamazaki, K. (2013, October). History matching for exploring and
 1352 reducing climate model parameter space using observations and a large
 1353 perturbed physics ensemble. *Climate Dynamics*, *41*(7-8), 1703–1729.
 1354 (WOS:000324812200002) doi: 10.1007/s00382-013-1896-4

- 1355 Williamson, D., & Volodina, V. (2020). Exeteruq mogp an r interface to performing
1356 uq with mop emulator. *Documentation*. Retrieved from [https://bayesexeter](https://bayesexeter.github.io/ExeterUQ_MOGP/)
1357 [.github.io/ExeterUQ_MOGP/](https://bayesexeter.github.io/ExeterUQ_MOGP/)
- 1358 Wurps, H., Steinfeld, G., & Heinz, S. (2020, March). Grid-Resolution Requirements
1359 for Large-Eddy Simulations of the Atmospheric Boundary Layer. *Boundary-*
1360 *Layer Meteorology*, *175*, 179–201. doi: 10.1007/s10546-020-00504-1
- 1361 Zhang, M., Somerville, R. C. J., & Xie, S. (2016, February). The SCM concept and
1362 creation of ARM forcing datasets. *Meteorological Monographs*, *57*, 24.1–24.12.
1363 doi: 10.1175/AMSMONOGRAPHIS-D-15-0040.1
- 1364 Zhang, Y., Klein, S. A., Fan, J., Chandra, A. S., Kollias, P., Xie, S., & Tang, S.
1365 (2017, October). Large-Eddy Simulation of Shallow Cumulus over Land: A
1366 Composite Case Based on ARM Long-Term Observations at Its Southern
1367 Great Plains Site. *Journal of the Atmospheric Sciences*, *74*(10), 3229–3251.
1368 doi: 10.1175/JAS-D-16-0317.1