

# Investigating cloud–radiation error compensations stemming from climate model calibration

M. Coulon-Decorzens<sup>1</sup>, F. Hourdin<sup>1</sup>, N. Villefranque<sup>2</sup>

<sup>1</sup>Laboratoire de Meteorologie Dynamique, Sorbonne Université/IPSL/CNRS, Paris, France

<sup>2</sup>Centre National de Recherches Météorologiques, Météo-France, CNRS, Toulouse, France

## Key Points:

- Ensembles of perturbed-parameter simulations of single column models are used to investigate cloud–radiation compensating errors
- We find that structural errors in radiative transfer models due to lack of 3D effects are compensated by overestimated cloud fractions
- Radiative transfer models with and without 3D effects should not be used with the same overlap and heterogeneity parameters

---

Corresponding author: M. Coulon-Decorzens, [maelle.coulon-decorzens@lmd.ipsl.fr](mailto:maelle.coulon-decorzens@lmd.ipsl.fr)

## Abstract

Compensating errors are an obstacle to the development of climate models. We wonder if systematic errors in simulated cloud properties might result from forced error compensation due to targeting top-of-atmosphere radiative fluxes in the tuning process while using an inaccurate radiative transfer parameterization. Here, we investigate parametric and structural errors in state-of-the-art radiative transfer models and ask if these errors might be compensated by errors in cloud properties, to what extent, and in which circumstances. Convection and cloud parameters of two versions of a Single-Column version of a climate Model (SCM), with and without a parameterization of 3D radiative effects of clouds, are tuned targeting reference solar fluxes calculated on Large-Eddy Simulation cloud fields. When 3D effects are neglected, reasonable fluxes are obtained only at the expense of overestimated cloud fractions, partly compensating for underestimated cloud reflectivity at low sun. Targeting mean solar-angle fluxes instead of detailed ones removes this mechanism entirely.

## Plain Language Summary

Global climate models are used to understand the climate system and anticipate the consequences of global warming. They are based on a set of equations that represent the various aspects of the system. Among these, cloud physics and radiative transfer (light propagation) play a major role. The equations, called parameterizations, are based on our understanding of the physics, complexity–accuracy compromises, and include parameters that need to be calibrated against observations. Satellite-observed radiative fluxes are often targeted in this calibration process, with the risk of obtaining the right fluxes (close to the observed ones) for the wrong reasons, and in particular for the wrong cloud properties. In this study, we show that the radiative transfer models currently used in climate models present structural errors that could, through the calibration process, be balanced by introducing additional errors in cloud properties. This implies that persistent systematic biases in clouds simulated by global models could be partly due to structural errors in the radiative transfer model, rather than to inadequate modeling of cloud physics. We use a History Matching approach based on machine learning to shed light on these questions and study error compensations in detail, albeit in an idealized framework.

## 1 Introduction

General circulation models (GCMs) used for climate projections are, like any model, imperfect representations of the climate system. Their behavior depends on free parameters that need to be adjusted, which is achieved through calibration. When calibrating numerical models as complex as GCMs, it is very difficult, if not impossible, to remove compensating errors. Reducing them is a motivation to develop more physically robust parameterizations. Finding ways to detect, characterize, and understand sources of compensating errors remains a major challenge in climate modeling, one that we hope to address more effectively thanks to increased computing power and machine learning algorithms. This issue is crucial for the reliability of climate change projections.

Hourdin et al. (2017) report that a common practice in calibrating (tuning) climate models is to target observed top-of-atmosphere (TOA) radiative fluxes by adjusting parameters associated with the most uncertain processes controlling these fluxes: those related to clouds. In doing so, accurate TOA fluxes are often obtained at the expense of cloud-related compensating errors: between cloud properties and e.g. surface albedo or jet position (Hourdin et al., 2013), between low-, middle- and high-level clouds (Webb, Senior, Bony, & Morcrette, 2001; Nam, Bony, Dufresne, & Chepfer, 2012) or even between physical, optical, and radiative properties of a given cloud regime (Konsta et al., 2022).

63 To understand these errors, and in particular those involving cloud physics and ra-  
 64 diative transfer (RT), it is necessary to disentangle model errors stemming from each pa-  
 65 rameterization.

66 An important step was taken in this direction in the recent years, with the sem-  
 67 inal work of Couvreux et al. (2021) and Hourdin et al. (2021), who proposed a tuning  
 68 strategy relying on a hierarchy of model configurations, from Single Column Model (SCM)  
 69 to GCMs. In those studies, SCMs are tuned to ensure that boundary-layer clouds are  
 70 well simulated, using a few representative cases for which Large Eddy Simulations (LES,  
 71 i. e. 3D simulations with resolution of a few tens of meters on domains of a few tens of  
 72 km) are used as a reference. Following the History Matching approach proposed by Williamson  
 73 et al. (2013), this tuning of the SCM does not produce one optimal parameter set for bound-  
 74 ary layer parameterizations, but rules out parameter values that lead to unsatisfyingly  
 75 large errors in the simulated boundary layers. Then, the GCM is tuned in the remain-  
 76 ing “Not Ruled-Out Yet” (NROY) parameter space. This limits error compensations be-  
 77 tween boundary layer physics and the other components of the model.

78 Here, we go one step further by investigating compensation errors that may occur  
 79 (i) *inside* one parameterisation, and (ii) between the parameterisations of the SCM, in  
 80 cumulus clouds situations, and focus on their solar radiative effect. In particular, we in-  
 81 vestigate (i) compensations that may arise between different aspects of cloud-geometry  
 82 inside a radiative transfer scheme and (ii) errors in convection and cloud schemes that  
 83 may arise from the need to compensate errors in the radiative transfer scheme when ra-  
 84 diative transfer is inaccurate.

85 Indeed, until recently most radiative transfer schemes implemented in GCMs ne-  
 86 glected or misrepresented important aspects of cloud geometry such as non-maximally  
 87 overlapped cloud fractions inside low-cloud layers, internal inhomogeneity of cloud wa-  
 88 ter, and 3D radiative effects associated with horizontal transport of light through cloud  
 89 sides, all of which affect cloud reflectivity. As detailed geometrical properties of clouds  
 90 vary in space and time, and might evolve in a warmer climat, it is important that their  
 91 individual effects, as well as their non-linear interactions, be well understood and cap-  
 92 tured in radiative transfer schemes. However, Villefranque et al. (2021) has shown that  
 93 the right reflectivity could be obtained for quite different sets of cloud-geometry param-  
 94 eters for cumulus clouds. Thus we ask: how to ensure that reflectivity is accurate for the  
 95 right reasons, i.e., for the right cloud-geometry parameters? On the other hand, we ac-  
 96 knowledge that cloud reflectivity errors will probably remain in the next generation of  
 97 GCMs; yet errors on cloud reflectivity will lead to errors on radiative fluxes even if ver-  
 98 tical profiles of cloud fraction and cloud water content are perfectly simulated by the bound-  
 99 ary layer parameterisations. Thus we ask: to which extent might reflectivity errors be  
 100 compensated by cloud property errors, and thus lead to the right fluxes for the wrong  
 101 reasons?

102 To answer these questions, we design cheap, idealized, SCM tuning experiments,  
 103 extending on previous work of Couvreux et al. (2021); Hourdin et al. (2021); Villefranque  
 104 et al. (2021). The specific radiative scheme and SCM that will be used in our study are  
 105 described in Section 2, as well as reference data used as target for the calibration, and  
 106 details on the History Matching tuning protocol. Section 3 is dedicated to the question  
 107 of error compensations inside the radiative scheme, and Section 4 to the question of radiation-  
 108 cloud error compensations. In Section 5 we discuss our results and their implications,  
 109 as well as the methodological choices on which they are based.

110 *[[ci dessous une liste des grev qui ont sauté vu les modifs sur l'intro]]*

- 111 • *[rev2 2)Some references are not cited correctly. Mais je ne vois pas l'erreur ?????]*
- 112 • *[rev2 3) sur les problèmes d'acronymes]*

- 113 • [rev2 1) "Words such as "Spartacus," "Tripleclouds," and "History Matching" should
- 114 preferably be introduced after their first occurrences]
- 115 • [Pour être plus explicite sur la première occurrence de History Matching, même si
- 116 c'est la deuxième ici. Remarque 1 rev2]
- 117 • [rev2 3) The full names and their abriavations are miused throughout the manuscript]
- 118 • [rev3 : This sentence is confusing. Would it be possible to add a few words more
- 119 precise then "more appropriate ecRad parameters"? I was confused by the words
- 120 "(instead of best Spartacus)". It seems trivial to me that using the parameters tuned
- 121 with Spartacus would not work well with Tripleclouds.] [MCD : D'ailleurs, c'est
- 122 surtout le fait de viser l'angle moyen qui enlève la surestimation systématique de
- 123 la clous cover nan ?] (la phrase était : Then, using more appropriate ecRad pa-
- 124 rameters for the 1D-RT scheme (instead of parameters values tuned for 3D-RT),
- 125 and targeting one SZA-averaged flux (instead of three SZA-dependent fluxes), we
- 126 show that both clouds and radiative fluxes can be reasonably well simulated in
- 127 the experiments using the 1D-RT scheme.)

## 128 2 Data, models and tools

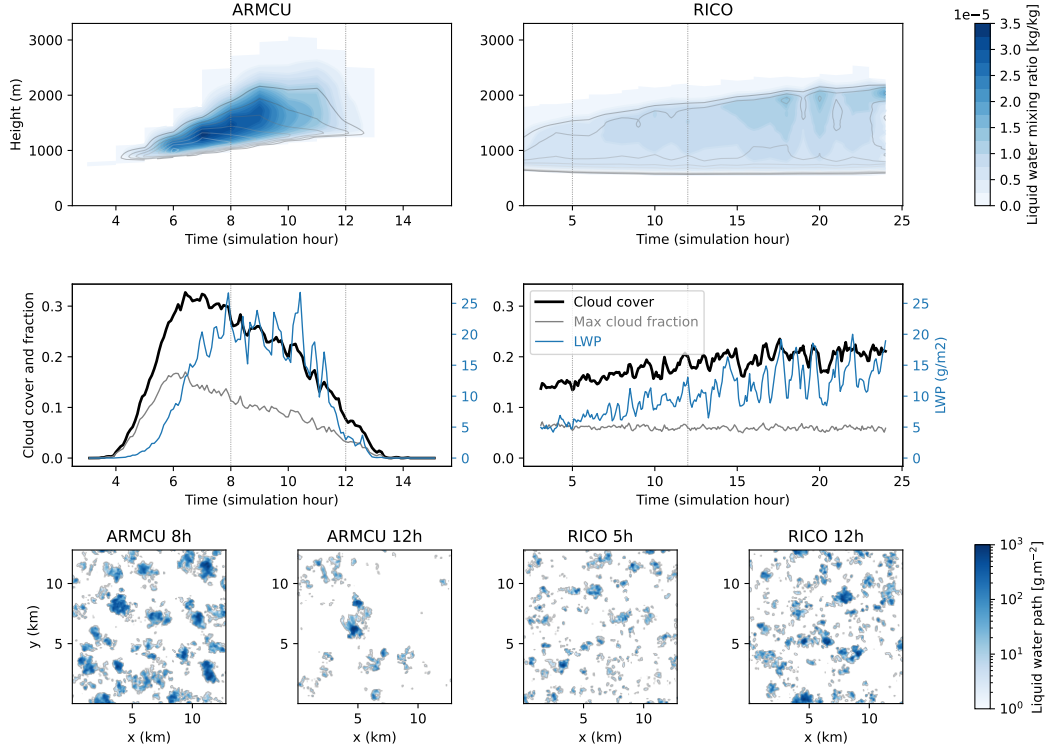
129 To run tuning experiments, three ingredients are needed: reference data (simulated  
130 or observed) that will serve as targets for the tuning; a model, which includes free pa-  
131 rameters that will be adjusted in the tuning process; a tuning tool (methodology and  
132 software) that will make the protocol objective and automatic.

133 In this study, two different kinds of tuning experiments are explored. They both  
134 use the same tuning tool, which is described in Section 2.4, and the same reference data:  
135 radiative fluxes computed using a 3D Monte Carlo code, run on 3D cloud fields output  
136 from LES, described in Section 2.1. All experiments are made using idealized cumulus  
137 cases, following Villefranque et al. (2019).

138 The difference between the two types of experiments is the model that is tuned:  
139 In Section 3, cloud-geometry parameters of the radiative transfer schemes of ecRad (Hogan  
140 & Bozzo, 2018) are tuned in a "perfect cloud" setup (ecRad is run on vertical profiles  
141 of cloud fraction and water content taken from horizontally-averaged 3D reference LES  
142 fields). The radiative transfer scheme and parameters are described in Section 2.2. In  
143 Section 4, cloud and boundary-layer parameterizations of the SCM version of LMDZ are  
144 tuned, targeting radiative fluxes. The SCM and its parameters are described in Section 2.3  
145 and Supporting Information Table S1.

### 146 2.1 Reference data

147 Two typical cumulus cloud cases are used from the set of idealized cases that are  
148 distributed in a standardized format by Dephy (Développement et Evaluation PHYsiques  
149 des modèles atmosphériques, Dephy (2020)). The ARMCU case (Brown et al., 2002) is  
150 typical of the development of boundary-layer clouds over continent during the day, while  
151 the RICO case (vanZanten et al., 2011) is typical of trade-wind cumulus developing over  
152 a stationary ocean. LES of these two cases are run with the Meso-NH model (Lafore  
153 et al., 1998; Lac et al., 2018) at 25 metres horizontal and vertical resolutions on a  $12.8$   
154  $\times 12.8 \times 4 \text{ km}^3$  domain. Large scale dynamics, radiative cooling and surface conditions  
155 are imposed throughout the simulation for each case. The effects of large-scale advec-  
156 tion and radiative cooling are represented by prescribed source terms in the heat and  
157 moisture evolution equations, applied in each column of the domain. These sources/sinks  
158 of heat and moisture are functions of height and time and replace an explicit resolution  
159 of the large-scale dynamics and of radiation. In the ARMCu case, surface fluxes are pre-  
160 scribed as a function of time, following a typical diurnal cycle of turbulent fluxes over  
161 land, while in RICO, the sea surface temperature is set to a constant during the whole  
162 simulation and turbulent fluxes are calculated by the model. [Demande de précision de



**Figure 1.** Reference simulations of the ARMCU (left) and RICO (right) cases used in this study. All variables are computed from LES 3D cloud fields. a-b: time evolution of vertical profiles of cloud fraction (contour lines, 0 to 0.16 every 0.02 in ARMCU, 0 to 0.07 every 0.01 in RICO) and cloud water mixing ratio (blue shading). c-d: time evolution of total (vertically integrated) cloud cover, maximum cloud fraction and domain-mean liquid water path. e-h: liquid water path maps at the two hours retained for tuning (ARMCU 8th and 12th hours, RICO 5th and 12th hours, identified by vertical lines in a-d plots).

163 *la part de reviewer 1 sur cette phrase : q2.2 Please elaborate on this. Naj: j'ai rajouté*  
 164 *des détails, dites moi si vous pensez que ça suffit ou pas ?].* These simulations provide  
 165 reference values for the thermodynamic and cloud variables, and their uncertainties are  
 166 quantified running sensitivity experiments to numerical and physics options as described  
 167 in Couvreur et al. (2021). The reference simulations are presented in Figure 1. Four scenes  
 168 from these two simulations will be used as constraints in the tuning experiments. They  
 169 are documented in Table 1.

170 Reference solar fluxes are computed using a 3D Monte Carlo (MC) code run on 3D  
 171 cloud fields extracted every hour from the LES, as described in Villefranque et al. (2019).  
 172 3D fields of liquid water content are taken from the LES and cloud-droplet effective radius  
 173 is homogeneously set to  $10 \mu\text{m}$ . Cloud optical properties are obtained from Mie theory.  
 174 Gas optical properties are calculated using the k-distribution model RRTMG-IFS  
 175 included in the ecRad radiation scheme (Hogan & Bozzo, 2018), for temperature, pressure  
 176 and humidity profiles corresponding to the LES horizontal mean below 4 km, and to  
 177 Standard Mid-Latitude Summer profile above. The solar constant is set to  $1368 \text{ W}\cdot\text{m}^{-2}$   
 178 and the surface albedo to 0.08. For each 3D cloud field, additional MC calculations are  
 179 made under the Independent Columns Approximation (ICA) by solving radiative transfer  
 180 independently in each column of the 3D LES cloud field and taking the average flux.  
 181 Assuming independent columns is the same as removing 3D radiative effects from the

**Table 1.** Cloud properties in four scenes extracted from ARMCU and RICO LES.

Case, hour	ARMCU, 8h	ARMCU, 12h	RICO, 5h	RICO, 12h
Cloud cover (%)	26.07	6.99	13.43	20.12
Maximum cloud fraction (%)	11.11	3.09	5.45	5.96
Cloud layer depth (km)	1.625	1.20	1.25	1.675
Cloud cov. / max. cloud frac.	2.346	2.26	2.462	3.379

182 calculation; the resulting flux is hence called 1D MC flux. Taking the difference between  
 183 3D and 1D MC fluxes provides an estimate of “3D radiative effects”.

184 [MCD : harmoniser ICA et 1D dans le texte]

## 185 2.2 Radiation parameterization

186 The radiative models under investigation in this study are the **Tripleclouds** and  
 187 **Spartacus** solvers implemented in ecRad, the radiative transfer model developed at Eu-  
 188 ropean Centre for Medium-Range Weather Forecasts (Hogan & Bozzo, 2018). ecRad pro-  
 189 vides a flexible interface that allows users to configure various aspects of the radiation  
 190 model. Cloud droplet effective radius, gas optics, clear-sky profiles (gas concentrations,  
 191 temperature and pressure) and radiative boundary conditions are set as in the MC sim-  
 192 ulations so that they are excluded from causes of possible differences between param-  
 193 eterized and reference fluxes. In “perfect clouds” experiments, input liquid water con-  
 194 tent and cloud fraction profiles are taken from horizontally averaged LES 3D fields and  
 195 are hence also excluded from potential causes of differences between parameterized and  
 196 reference fluxes. In SCM experiments, liquid water content and cloud fraction profiles  
 197 input to ecRad are taken from SCM simulation outputs.

198 Cloud optics are interpolated from a Mie look-up table provided with ecRad, sim-  
 199 ilar to, but slightly different from, the one used in the Monte Carlo simulations. The dif-  
 200 ference between optical properties taken from the two tables are typically less than a change  
 201 of 1  $\mu\text{m}$  in the effective radius of cloud droplets (not shown). [rev3 2.3]

202 The RT models at the heart of **Tripleclouds** and **Spartacus** are modified ver-  
 203 sions of the two-stream model (Meador & Weaver, 1980), which directly integrate the  
 204 effects of cloud geometry on radiation transport through assumptions on vertical over-  
 205 lap, horizontal heterogeneity and, in **Spartacus** only, cloud size.

206 In our configurations of both **Tripleclouds** and **Spartacus**, vertical overlap is rep-  
 207 resented using the exponential-random model parameterized by its decorrelation length  
 208  $\ell$  (Hogan & Illingworth, 2000). A two-region cloud representation (the **Tripleclouds**  
 209 model of Shonk and Hogan (2008)) is used to account for in-cloud water sub-grid het-  
 210 erogeneity, whereby layer-wise optical depths in thin-cloud and thick-cloud regions are  
 211 calculated according to the fractional standard deviation (*FSD*) parameter. In the **Tripleclouds**  
 212 solver, no 3D effects are taken into account, whereas in **Spartacus** (Hogan, Schäfer, Klinger,  
 213 Chiu, & Mayer, 2016; Schäfer, Hogan, Klinger, Chiu, & Mayer, 2016; Hogan, Fielding,  
 214 Barker, Villefranche, & Schäfer, 2019), intensity of 3D effects is proportional to cloud-  
 215 side perimeter length (Hogan & Shonk, 2013), itself a function of cloud fraction and cloud  
 216 effective scale ( $C_s$ ). In the rest of this article, the **Tripleclouds** solver is referred to as  
 217 1D-RT and the **Spartacus** solver as 3D-RT, referring to the presence or absence of 3D  
 218 radiative effects.

219

### 2.3 LMDZ Single-Column Model

220

221

222

223

224

225

LMDZ-6A (Hourdin et al., 2020) is the atmospheric component of the IPSL-6A General Circulation Model, which participated in the sixth phase of the Coupled Model Intercomparison Project (CMIP6). Here, its single-column version is used with a refined 95-level grid as in Hourdin et al. (2019, 2021) to simulate ARMCU and RICO cases. The same large-scale dynamics, radiative trends and surface conditions are imposed as in the LES so that physical parameterizations are the only active part of the model.

226

227

228

229

230

231

232

233

234

235

236

237

238

239

240

241

242

243

More specifically, the parameterizations that are active here are the boundary-layer transport and cloud schemes. The parameterization of vertical sub-grid transport is based on an Eddy-Diffusivity and Mass-Flux approach. The Eddy-Diffusivity model parameterizes the effects of small-scale turbulence on the mean state using the Turbulent Kinetic Energy prognostic equation formulated by Yamada (1983) with a 1.5-order closure. The Mass-Flux model parameterizes the effects of organized convective cells or rolls on the mean state using an effective thermal plume model. The plume transports air and state variables from the surface to the boundary-layer top. Exchanges with the environment are modeled through lateral entrainment and detrainment formulations (Hourdin et al., 2019). Water condensate and cloud fraction profiles are computed using a bi-Gaussian probability density function of the saturation deficit, with one mode accounting for saturation deficit in the thermal plume and one mode in the environment (Jam, Hourdin, Rio, & Couvreur, 2013). This combination of Eddy-Diffusivity-Mass-Flux scheme with a bi-Gaussian cloud scheme provides a unified framework that has been shown to accurately represent both dry and cloudy convective boundary layers with cloud regimes ranging from cumulus to stratocumulus (Hourdin et al., 2019). The conversion of cloud water into precipitation and the evaporation of precipitation are detailed in Madeleine et al. (2020).

244

245

246

247

248

249

250

251

252

253

254

As for radiation, ecRad was recently implemented in LMDZ and will be part of the forthcoming versions of the GCM. However, in this study, it is the offline version of ecRad that is run on the SCM output profiles. It calculates the radiative fluxes associated with the SCM cloud profiles and ecRad cloud-geometry parameters. The offline version was used to make more robust comparison, focused on cloud radiative effects. In particular it allows to use the same clear-sky profiles as in the MC simulation. This was done in practice by replacing the clear-sky profiles in the SCM outputs by the profiles of the LES averaged horizontally. Note that in the two cumulus cases that are used as constraints, surface fluxes and radiative cooling are prescribed as forcings: radiation is not interactive in the simulations and thus does not affect the clouds (consistently in the LES and SCM).

255

### 2.4 The High-Tune:Explorer tuning tool

256

257

258

259

260

261

262

High-Tune:Explorer is a tuning tool based on History Matching with iterative refocusing (Vernon, Goldstein, & Bower, 2010; Williamson et al., 2013). It aims at finding the subspace of model free parameters that match a set of constraints. To that end, the parameter space hypercube,  $[\lambda_{min}^1, \lambda_{max}^1] \times [\lambda_{min}^2, \lambda_{max}^2] \times \dots \times [\lambda_{min}^N, \lambda_{max}^N]$ , with  $\lambda^1, \dots, \lambda^N$  the  $N$  free parameters to tune, is iteratively reduced by ruling out parameter vectors for which the model's predictions, for a set of user-defined metrics, do not match reference values within the range of user-defined tolerance to error.

263

264

265

To accelerate the exploration of the hypercube, emulators based on Gaussian Processes are built for each metric. These emulators are trained on metrics computed from an ensemble of model runs (with typically  $10 \times N$  members), to then provide rapid pre-

266 dictions of metric values for huge sets of free parameter vectors <sup>1</sup>. Prediction uncertainty  
 267 is added to tolerance-to-error in order to avoid ruling out parameter vectors that might  
 268 be far from the target because of bad emulators but actually acceptable configurations  
 269 of the model. To this end, implausibility is defined for any parameter vector  $\boldsymbol{\lambda}$  as

$$I(\boldsymbol{\lambda}) = \max \left\{ \frac{|r_1 - \mu_1(\boldsymbol{\lambda})|}{\sqrt{\sigma_1^2(\boldsymbol{\lambda}) + T_1^2}}, \dots, \frac{|r_p - \mu_p(\boldsymbol{\lambda})|}{\sqrt{\sigma_p^2(\boldsymbol{\lambda}) + T_p^2}} \right\}, \quad (1)$$

270 with  $r_i$  the reference (target) value and  $T_i$  its tolerance to error.

271 The parameter vector  $\boldsymbol{\lambda}$  is ruled out if its implausibility  $I(\boldsymbol{\lambda})$  is greater than an ar-  
 272 bitrary value  $\Gamma$ , which represents the size of the confidence interval (reference  $\pm\Gamma$  times  
 273 uncertainty), typically chosen between 2 and 3. At the end of each iteration, the new Not-  
 274 Ruled-Out-Yet (NROY) space of parameters is determined using this implausibility con-  
 275 dition. The next iteration starts by sampling a set of parameter vectors in the NROY  
 276 space of the previous iteration. [*rev1 q2.3 : Have the samples that were generated dur-*  
 277 *ing previous iterations (but still within the NROY space) been reused when training the*  
 278 *emulator?*] Then a new ensemble is run, metrics are evaluated, emulators are built, etc.  
 279 In practice, the parameter vectors used to run the first ensemble of simulation are sam-  
 280 pled randomly in the hypercube using a latin hypercube sampling algorithm. The en-  
 281 semble of next waves used the  $10 \times N$  first members that satisfies the implausibility con-  
 282 dition, among another bigger latin hypercube sampling of typically  $10^5$  samples. [*Pour*  
 283 *répondre à la question 2.4 de rev1*].

284 As the iterative process progresses, the NROY space narrows down, mostly because  
 285 emulators uncertainty decreases, which is due to denser information being collected for  
 286 training (same amount of points in a smaller NROY space). The tuning experiment is  
 287 considered to have strictly converged when emulator uncertainties  $\sigma$  are significantly smaller  
 288 than tolerances to error  $T$  for every metrics. In that case, the final NROY space is ex-  
 289 actly the subspace of free parameters that matches the user-defined constraints and em-  
 290 ulators can be considered reliable models for the metrics (inside the final NROY space).

291 In practice, emulator uncertainties rarely fall one order of magnitude under toler-  
 292 ances to error for all metrics and hence the experiment rarely strictly converges. We there-  
 293 fore find it useful to define weak convergence: The experiment is considered to have weakly  
 294 converged when adding new iterations does no longer significantly reduce NROY space.  
 295 In that case, the final NROY space is larger than the sought parameter subspace.

296 Importantly, during the tuning experiment, many simulations will have been run,  
 297 constituting a Perturbed-Parameter Ensemble (PPE) [*MCD : J'ai déjà introduit la sig-*  
 298 *nification de PPE dans l'intro : je la remet ici ou pas ? Naj : il y est plus dans l'intro*  
 299 *je crois !*]. For each metric  $i$  and each simulation  $j$ , its distance to reference value can  
 300 be calculated as  $\frac{|r_i - f_i^j|}{T_i}$ , where  $f_i^j = f_i(\boldsymbol{\lambda}^j)$  is the actual model output for metric  $i$  and  
 301 parameter vector  $\boldsymbol{\lambda}^j$ . The score of a simulation  $j$  is then defined as its worst-metric score:

$$S(\boldsymbol{\lambda}^j) = \max \left\{ \frac{|r_1 - f_1(\boldsymbol{\lambda}^j)|}{T_1}, \dots, \frac{|r_p - f_p(\boldsymbol{\lambda}^j)|}{T_p} \right\}, \quad (2)$$

<sup>1</sup> Training an emulator means characterizing the law of the Gaussian Process  $G \sim \mathcal{GP}(m, k(\cdot, \cdot))$ , that is, specifying its expectation function  $m$  and variance kernel  $k(\cdot, \cdot)$ . This is done by fitting standard analytic functions to the simulated data for metric  $i$ ,  $\{f_i^j\}_{j \in 1, \dots, 10N}$  obtained from model runs at  $10N$  points  $\{\boldsymbol{\lambda}^j\}_{j \in 1, \dots, 10N}$  sampled in parameter space. Then, for any new  $\boldsymbol{\lambda}^*$  in parameter space, the  $i$ -th emulator provides a prediction of the  $i$ -th metric as the expectation  $\mu_i$  and standard deviation  $\sigma_i$  of the conditional random variable  $G_i^* = G_i(\boldsymbol{\lambda}^*) | \{G_i(\boldsymbol{\lambda}^j) = f_i^j, \forall j \in [1, 10N]\}$ ; and  $G_i^* \sim \mathcal{N}(\mu_i, \sigma_i)$ .

where emulator predictions and uncertainties no longer appear, contrary to the implausibility defined in Equation (1). Keeping only simulations whose scores are under a given threshold provides, independently from the emulators convergence status, constrained PPEs that can be studied to learn about the quality of the model. Of course, if the emulators convergence is bad, the NROY space might remain too large and model configurations sample might not be dense enough, so low-probability, acceptable model behaviours might be absent from the PPEs.

### 3 The right cloud reflectivity for the right cloud-geometry parameters?

In this section, HighTune:Explorer is used to explore the behaviour of the 1D- and 3D-RT solvers in a “perfect-cloud” setup. Their cloud-geometry parameters are varied to produce PPEs that are then analysed. These parameters are the overlap decorrelation length  $\ell$ , the fractional standard deviation of in-cloud water  $FSD$ , and the cloud effective scale  $C_s$  (only active in the 3D-RT model). As the focus is on internal RT compensating error related to cloud geometry, flux differences that might be related to physical cloud properties are removed by using mean LES cloud fraction and liquid water content profiles as inputs to the 1D- and 3D-RT solvers. These profiles were horizontally-averaged over the 3D fields that were used as inputs to the Monte Carlo reference RT simulations.

Effects of cloud geometry on cloud reflectivity have been studied for a long time (see e.g., McKee and Cox (1974); Barker, Stephens, and Fu (1999); Várnai and Davies (1999); Shonk, Hogan, Edwards, and Mace (2010); Hogan et al. (2019); Villefranque et al. (2021)). *[rev2 2) Some references are not cited correctly. Mais je ne vois pas l'erreur ????? Naj : je pense que c'est la liste de noms au lieu du "et al." mais c'est le format de JAMES on y peut rien]* In particular for solar reflectivity of cumulus clouds, it is well-known that:

- the widely used maximum overlap assumption tends to underestimate total (i.e. vertically integrated) cloud cover, and consequently, cloud reflectivity; *[MCD : on a précisé qu'on parlait du SW ici ? Naj : maintenant oui !]*
- neglecting in-cloud optical-depth inhomogeneity systematically leads to overly reflective clouds;
- assuming infinite cloud size means neglecting 3D radiative effects, which produces clouds that are too bright by around  $10\text{W}\cdot\text{m}^{-2}$  when the sun is at zenith, and too dim by around  $10\text{W}\cdot\text{m}^{-2}$  when the sun is low on the horizon. Note that since 3D effects shift from positive to negative values depending on Solar Zenith Angle (SZA), *[rev3 3) sur les problèmes d'acronymes Naj : ok c'est bien la première occurrence ici pour l'instant]* they are null at around  $40\text{-}50^\circ$ , and partially cancel out over a day. *[rev3 2.4 : Does it account for  $10\text{W}/\text{m}^2$  for high and low sun? J'ai pas compris ce qui allait pas dans la phrase Naj : la phrase a changé, est ce que c'est plus clair maintenant ? (J'avais pas compris non plus la question)]* *[[je déplace ça de la sec 3.2 à ici car pas vraiment un nouveau résultat - mais j'ai pas encore repris / vérifié ta phrase]* This can be explain by the competing effects between side leakage, that lead to less reflective clouds and dominate at low SZA, and side illumination that lead to more reflective clouds and dominate at high SZA. *[MCD : Une ref ? Je veux bien que vous vérifier les sens.]* *[Pour répondre à rev2 q5) Another interesting thing in section 3 is that the results in both spartacus and tripleclouds are close when the SZA is  $44^\circ$ . What is the underlying mechanism?]*

Recent RT models such as those implemented in ecRad Shonk and Hogan (2008); Hogan and Shonk (2013); Hogan et al. (2019); Pincus, Barker, and Morcrette (2003) include terms that represent effects of cloud geometry on radiative fluxes, and are able to simulate fairly

352 accurate cloud reflectivity when compared to reference 3D MC simulations. Yet, differ-  
 353 ent choices of cloud-geometry parameter values may result in the right cloud reflectiv-  
 354 ity. Here, we investigate these configurations, to determine if the right reflectivity might  
 355 be achieved for the wrong reasons, that is, for the wrong cloud geometry.

### 356 3.1 Experimental design

357 In all experiments, the three-dimensional parameter space that is explored is  $\Omega_{\ell, FSD, C_s} =$   
 358  $\square \times \square \times \square$  [MCD : c'est laissé vide pour mettre les bornes de chaque paramètre ?]. All  
 359 tuning experiments target solar reflected fluxes of the four scenes of ARMCU and RICO  
 360 cases (see Figure 1) for three distinct SZA,  $0^\circ$ ,  $44^\circ$  and  $77^\circ$ , following Villefranque et al.  
 361 (2021). This choice of distinct SZA values is made to have metrics that are sensitive to  
 362 3D radiative effects. [MCD : J'enleve l'histoire des signes opposés car c'est dit plus haut]

363 Reference radiative fluxes used as targets are 3D-MC fluxes for the 3D-RT tuning  
 364 experiments, and 1D-MC fluxes for the 1D-RT experiments. It is particularly significant  
 365 that we require the 1D-RT model to perform well in 1D and not in 3D: we do not force  
 366 it to compensate for a structural error that we know is present in the solver, namely the  
 367 absence of 3D radiative effects. Using 1D-MC references is a common practice in 1D-  
 368 RT model evaluation (see e.g. Barker et al. (1999), among many), though it is not nec-  
 369 essarily the strategy that would effectively be employed in the tuning of a GCM, where  
 370 "true" radiative fluxes might be targeted even at the expense of internal compensating  
 371 errors.

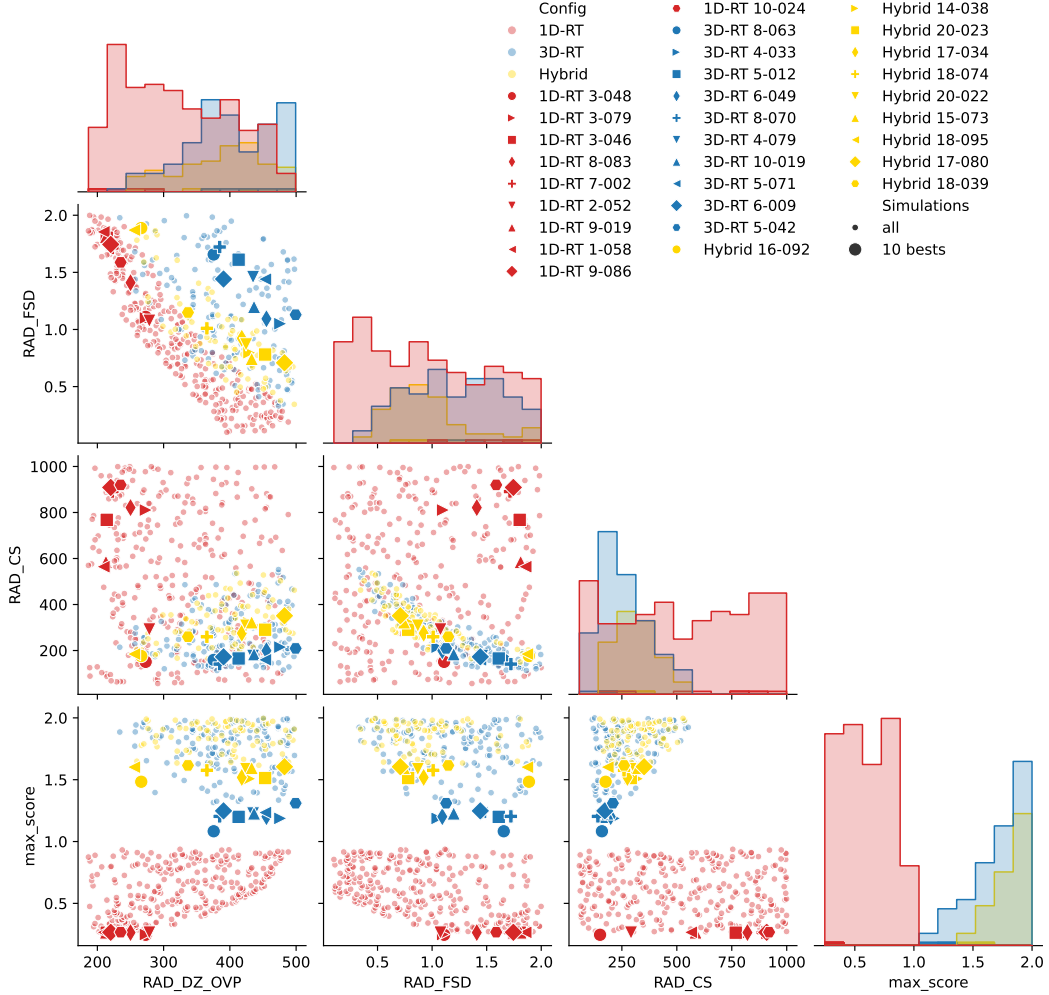
372 In the tuning experiments, parameter values that result in errors larger than three  
 373 times the tolerance to error  $T = 3.5 \text{ W}\cdot\text{m}^{-2}$  are rejected. This value of tolerance to er-  
 374 ror is an estimate of the **Spartacus** structural error on cumulus clouds, based upon the  
 375 work of Villefranque et al. (2021) (see Supporting Information Text S1).

376 For each experiment, ten waves (iterations) of history matching are produced, each  
 377 consisting in 101 simulations. [q2.4 reviewer1 : Which sampling algorithm was applied?  
 378 e.g., LHS? QMC? Naj : on l'a pas dit dans la section 2 ? Maelle Si ça a été rajouté] In  
 379 total, for each experiment, 1010 configurations of the solver (either **Tripleclouds** or  
 380 **Spartacus** depending on the considered experiment) are sampled and run. Among these  
 381 configurations, the 300 simulations with the lowest scores, i.e. the bests ones, are selected.  
 382 These ensembles constitute the PPEs under study in this section, named 1D-PPE for the  
 383 **Tripleclouds** tuning experiments and 3D-PPE for the **Spartacus** one.

### 384 3.2 Independent calibration of 1D- and 3D-RT solvers

385 Figure 2 shows the location of the 1D-PPE and 3D-PPE members in parameter  
 386 space. The first three rows and columns of the subplot matrix are relative to paramete-  
 387 r values, while the last row and column is relative to the score  $S(\boldsymbol{\lambda})$  of each configu-  
 388 ration, as per Equation (2). Since the  $C_s$  parameter is not used in the 1D-RT solver, the  
 389 red points corresponding to this configurations spread randomly along the  $C_s$  axis.

390 First, note that all configurations of the 1D-PPE have scores below 1, while all those  
 391 of 3D-PPE have scores above 1. Of course it does not mean that 1D-RT is better than  
 392 3D-RT in an absolute sense, since the scores of 1D-PPE are computed against 1D-MC  
 393 references, intentionally neglecting 3D effects. It just means that 1D-RT is closer to re-  
 394 ference 1D fluxes than 3D-RT is close to reference 3D fluxes. It might not be too surpris-  
 395 ing for a two-stream model to perform better at simulating 1D radiation than 3D, es-  
 396 pecially when considering the many years spent on the 1D problem (i.e., how to account  
 397 for vertical cloud structure and horizontal heterogeneity in vertical light propagation)  
 398 compared to the relatively short time dedicated to modelling horizontal transport within  
 399 the two-stream framework. In fact, when running the same tuning experiment with the  
 400 1D-RT targeting 3D-MC fluxes instead of the 1D-MC ones, tolerance to error had to be



**Figure 2.** Parameters of 300 best simulations for three “perfect-cloud” ecRad tuning experiments: 1D-RT *Tripleclouds* (red), 3D-RT *Spartacus* (blue) and Hybrid (gold). [MCD : faire en sorte de ne pas avoir besoin de la phrase qui suit] RAD\_DZ\_OVP is the vertical overlap decorrelation length parameter  $\ell$ , RAD\_FSD is the relative in-cloud inhomogeneity parameter  $FSD$ ,  $C_s$  is the cloud size parameter  $C_s$ . The max\_score is  $S(\lambda)$ . Each column corresponds to one parameter. The top subplot of each column shows histograms of the parameter values for the 300 best configurations, for each experiment. Then, each row corresponds to a second parameter, and subplots show the location of the configurations in the 2D space of the column  $\times$  row parameters.

401 increased from  $3.5 \text{ W}\cdot\text{m}^{-2}$  to  $9 \text{ W}\cdot\text{m}^{-2}$  to reach a non-empty NROY space, along with  
 402 an increase of the best score simulation (see ). Nevertheless, the NROY space of the two  
 403 experiment are very similar and simulated very similar radiative fluxes (not shown). The  
 404 difference of tolerance to error is in fact explain by the change of fluxes reference values,  
 405 and not by the change in fluxes simulation. This result mean that 1D-RT cloud geom-  
 406 etry parameter can't compensate for the lack of 3D radiative effects.

The best 1D-RT simulations, highlighted by large red symbols in Figure 2, exhibit a strong negative correlation between  $\ell$  and  $FSD$ . In 1D radiation, this is easy enough to understand: with a smaller decorrelation length, clouds are more randomly (or less maximally) overlapped, leading to a larger cloud cover and smaller cloud optical thickness in regions covered by clouds. The radiative effect of increasing cloud cover dominates that of decreasing cloud optical thickness, because reflectivity is an exponential function of optical thickness  $\tau$ ,

$$R \sim 1 - \exp(-\tau).$$

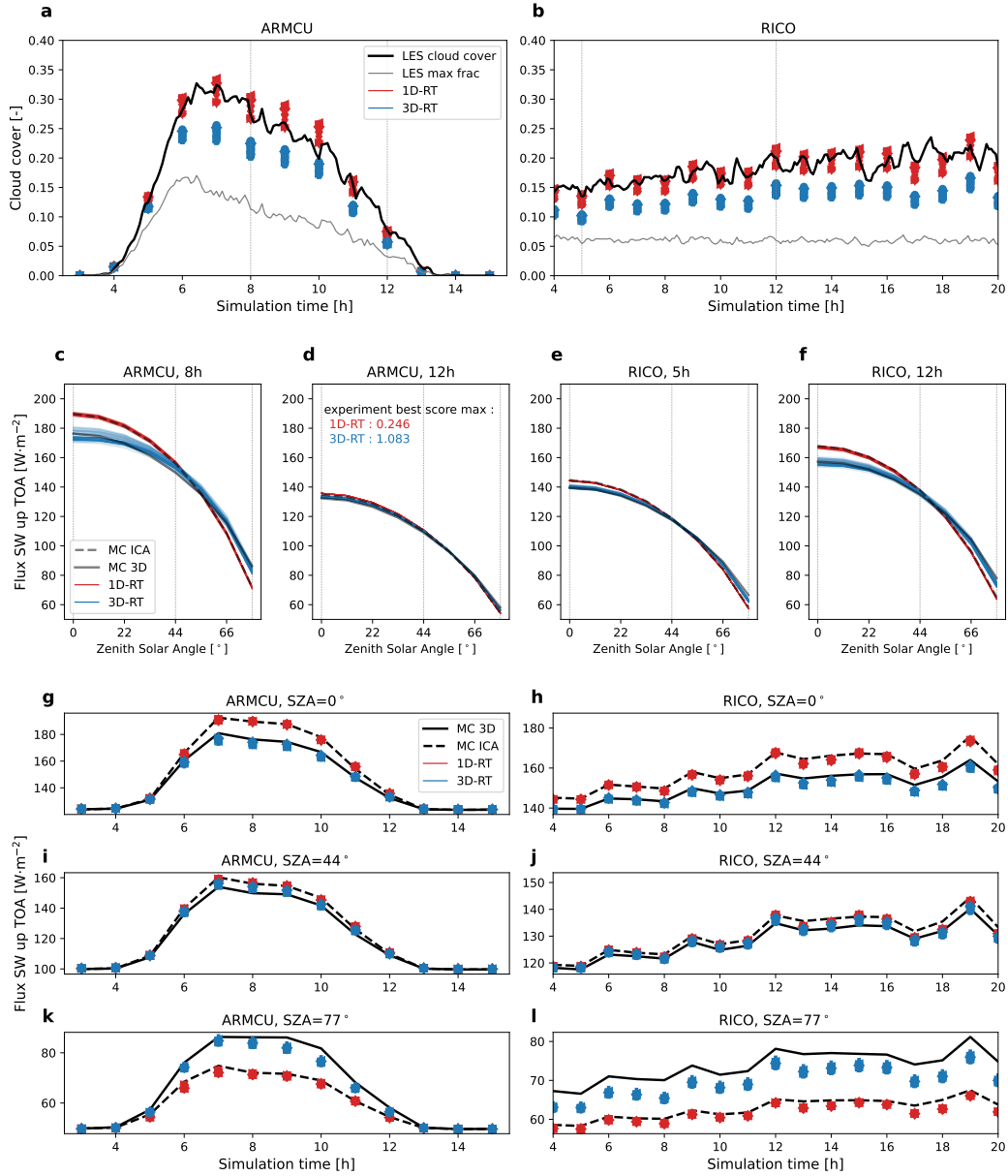
407 The optical thickness at a given location on the horizontal is the sum of the optical thick-  
 408 nesses of the cloudy layers above it. *[rev3 2.6 There is a part of this sentence missing. J'ai rajouté it mais pas sur que ça suffise]* Distributing the individual optical thicknesses  
 409 of the cloud layers horizontally on a wider area (large cloud cover with clouds of mod-  
 410 erate optical thickness) leads to a larger reflectivity than when cloudy layers overlap, be-  
 411 cause of the saturation of the exponential term.  
 412

413 The same reasoning can be applied to in-cloud water horizontal heterogeneity: be-  
 414 cause of the exponential function, horizontally homogeneous clouds are more reflective  
 415 than heterogeneous clouds that have the same horizontally-averaged optical depth. Thus,  
 416 decreasing  $\ell$  leads to more reflective clouds, and increasing  $FSD$  leads to less reflective  
 417 clouds. This explains why the tuning algorithm selects smaller values of the  $FSD$  to com-  
 418 pensate for larger values of  $\ell$  in order to simulate the right fluxes.

419 In the 3D-PPE (blue symbols on Figure 2), the relationship between  $\ell$  and  $FSD$   
 420 also seems to exist, but the 10 best configurations of this PPE show more dispersion than  
 421 the 10 best 1D-PPE ones. Our first guess was that horizontal transport, driven by the  
 422 third parameter  $C_s$ , modulates reflectivity in 3D-RT while it has no effect in 1D-RT, which  
 423 means that, in 3D-RT, a given value of  $FSD$  can be associated with various values of  
 424  $\ell$  and still produce the same reflectivity. For instance a smaller cloud cover could be com-  
 425 pensated by more intense 3D effects to maintain reflectivity at large SZA. However, 3D  
 426 effects tend to decrease reflectivity when the sun is close to zenith and increase it oth-  
 427 erwise; and  $C_s$  does not vary much between these 10 best configurations so this expla-  
 428 nation might not be sufficient.

429 All things considered, the internal dispersion of the best 3D-RT configurations might  
 430 be a detail compared to the fact that the 10 best configurations of each experiment are  
 431 located in distinct regions of the  $\ell \times FSD$  space: The best 3D-RT configurations are  
 432 located in the “large  $\ell$ , large  $FSD$ ” corner of the parameter space, where none of the  
 433 best 1D-RT simulations can be found. It implies that cloud geometry parameters yield-  
 434 ing strongly overlapped, highly heterogeneous clouds are not acceptable from a 1D-RT  
 435 perspective. Does this mean that reflectivity in the best 3D-RT configurations is right  
 436 for the wrong reasons?

437 To investigate this point further, let us look at cloud covers calculated by the exponential-  
 438 random overlap model, for the 10 decorrelation length values considered as “best” val-  
 439 ues in 1D-PPE compared to the 10 best values in the 3D-PPE. The upper panel of Fig-  
 440 ure 3 presents a comparison of reference cloud covers diagnosed in the 3D fields of ARMCU  
 441 and RICO LES simulations, with those calculated by the overlap scheme. Note that cov-  
 442 ers calculated in ecRad, when using 1D-RT or 3D-RT solvers, do not enter directly in  
 443 the flux equations. They are only computed as a diagnosis. Rather, the overlap model  
 444 is used locally, at each interface between layers, to distribute fluxes from cloudy and clear



**Figure 3.** [MCD : harmoniser MC ICA-1D] Cloud and radiation variables for the 10 best simulations of two perfect-cloud experiments (1D-RT with `Tripleclouds` in red, 3D-RT with `Spartacus` in blue), compared to reference. Left column: ARMCU. Right column: RICO. a-b: cloud cover (black) and maximum cloud fraction (gray) in LES fields, and cloud cover computed by eRad (color points) using input cloud fraction profiles calculated from LES cloud fields, and the exponential-random overlap model parameterized by various decorrelation lengths (RAD\_DZ\_OVP values in Figure 2). c-l: Fluxes simulated by Monte Carlo 3D (black full lines) and 1D (black dashed lines), 1D-RT with `Tripleclouds` (red) and 3D-RT with `Spartacus` (blue). c-f: upward TOA fluxes as a function of solar zenith angle, for the 4 cloud scenes used as constraints in the tuning experiments: ARMCU 8th (c) and 12th (d) hours, RICO 5th (e) and 12th (f) hours. For  $SZA=0^\circ$ ,  $44^\circ$  and  $77^\circ$ , fluxes plotted are directly those used as metrics. For the other one ( $SZA=11^\circ, 22^\circ, 33^\circ, 55^\circ, 66^\circ$ ), they are further computed following the same framework, for diagnosis only. [Pour répondre à rev1 q2.5 "it seems the experiments were run at more than three solar angles, but I didn't see any description about this configuration in the method section."] g-l: upward TOA fluxes, calculated at every hour of the two LES simulations, each time for three solar angles. g-h:  $SZA=0^\circ$ , i-j:  $SZA=44^\circ$ , k-l:  $SZA=77^\circ$ .

445 regions of a given layer, to cloudy and clear regions of adjacent layers. These plots show  
 446 that cloud covers corresponding to the best 1D-RT simulations are very close to those  
 447 of the LES. In comparison, cloud covers corresponding to the best 3D-RT simulations  
 448 are systematically smaller than those of the LES (in agreement with the fact that  $\ell$  values  
 449 are larger). This suggests that in these best 3D-RT configurations, reflectivity is right  
 450 for the wrong reasons, i.e., for clouds that are too maximally overlapped (we might say,  
 451 too compact). In this view, we conclude that compensating errors are at work between  
 452 overlap, inhomogeneity and 3D effects in the 3D-RT solver. This is an important result,  
 453 as it means that finding the best 3D-RT configuration and removing 3D effects (e.g., using  
 454 1D-RT with the same  $(\ell, FSD)$  values, or setting  $C_s$  to infinity in 3D-RT), will produce  
 455 wrong fluxes compared to the 1D MC reference.

456 A last important result shown in Figure 3(c-f) is that only the 3D-RT solver is able  
 457 to capture the angular dependency of cloud reflectivity. *[rev1 2.7 What is "the second line"?*  
 458 *Naj : en fait c'est parce que le bon mot est "row" ] [MCD : pas sur que ce soit bien dit*  
 459 *le "alone" dans la phrase, moi je compris qu'il a besoin de personne d'autres pour bien*  
 460 *faire la dépendance angulaire, alors qu'on veut plutôt dire que c'est le seul à pouvoir la*  
 461 *faire ? Ou bien j'ai mal compris ?]* The 1D-RT always underestimates (resp. overesti-  
 462 mates) cloud reflectivity at high SZA (resp. at low SZA) compared to 3D radiation. *[rev1*  
 463 *2.8 : Compared to zenith or compared to Spartacus? In addition, Tripleclouds was cal-*  
 464 *ibrated targeting 1D MC flux. Can we improve the angular dependency of Tripleclouds*  
 465 *if targeting 3D MC flux? Maëlle : remarque bien prise en compte dans ce paragraphe]*  
 466 *[[on l'a dit plus haut du coup j'enlèverai ça : Note that 1D-RT is equal to 3D-RT at mid-*  
 467 *dle SZA (around 44°) because 3D radiative effects generally cancel out around middle*  
 468 *SZA, as can be deduced from the 1D- and 3D-MC lines of Figure 3 (c-f). Maëlle : oui*  
 469 *bien sur]* This is not an artifact of the calibration experiment: tuning the 1D-RT scheme  
 470 against 3D-MC fluxes leads to the same result (see Figure S2 in Supplementary Infor-  
 471 mation).

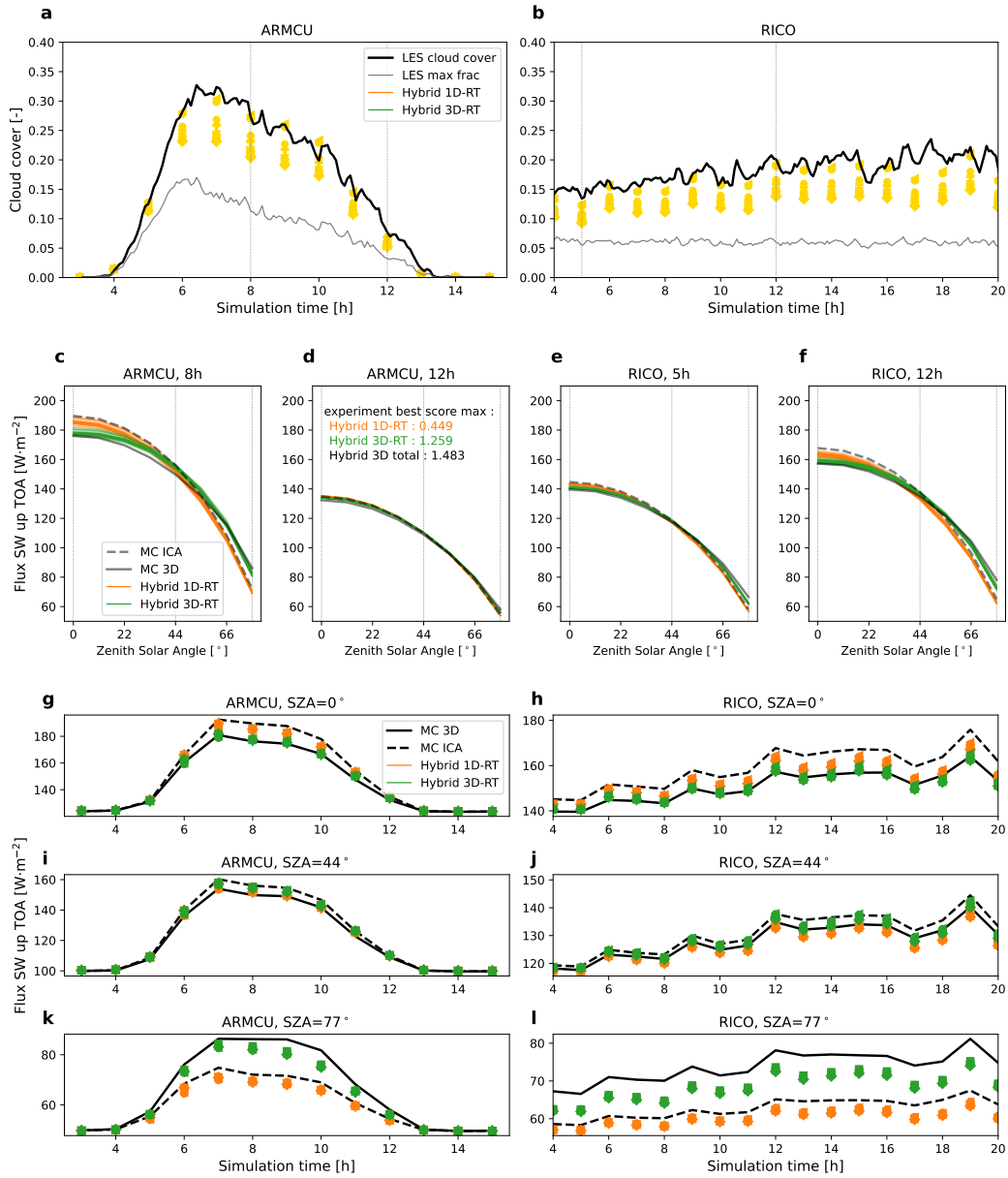
### 472 3.3 Obtaining the right reflectivity for the right cloud geometry

473 The best 3D-RT configurations were found to yield the right reflectivity for the wrong  
 474 cloud geometry, which we interpret as a sign of internal compensating errors. Our main  
 475 hypothesis here is that a "good" 3D-RT model should be able to simulate 3D fluxes as  
 476 well as 1D fluxes when 3D effects are removed from the simulation. In the following, we  
 477 look for configurations of the 3D-RT scheme that satisfy this requirement. *[Ce paragraphe*  
 478 *correspond a la remarque de rev3 : It is not clear to me why the authors try to find a pa-*  
 479 *ramete set that works well for the two different models simultaneously. Is it to reach the*  
 480 *objectives of section 4? Please clarify how at the beginning of this section. Et aussi a sa*  
 481 *remarque 2.8 Could you clarify here why you want to satisfy both solvers at the same time?]*

482 We thus design a third tuning experiment where  $(\ell, FSD)$  values that do not re-  
 483 sult in both accurate 3D fluxes, and accurate 1D fluxes when 3D effects are removed,  
 484 are ruled out from the acceptable 3D-RT parameter space. It is called the Hybrid ex-  
 485 periment, since it looks for  $(\ell, FSD, C_s)$  values that satisfy both solvers at the same time.

486 The 300 best Hybrid configurations are presented in Figure 2 (yellow marks). Their  
 487 scores are generally worse than for previous experiments (mostly, above 1.5). This is ex-  
 488 pected since the best 3D-RT configurations of Section 3.2 are bad choices for 1D-RT, and  
 489 conversely. However, a compromise was found within our tolerance to error.

490 The location of the best Hybrid configurations in the  $(\ell, FSD)$  space illustrates this  
 491 compromise: for a given choice of  $FSD$ , the 10 best simulations of the Hybrid experi-  
 492 ment correspond to  $\ell$  values between best 1D-RT and best 3D-RT values. Note that rel-  
 493 atively small values of  $FSD$  of around 0.75, associated with relatively large  $C_s$  of around  
 494 400 m, yield scores that are among the ten smallest ones, which was not the case in any  
 495 of the previous experiments.



**Figure 4.** *[rev1 2.10. Figure 4, What are the yellow marks? J'ai fait des changements plus loin pour répondre à cette question + il faudrait ajuster dans la légende du jaune pour Hybrid]* Cloud and radiation variables for the 10 best configurations of the Hybrid experiment, compared to reference. As in Figure 3 except that: in a-b, decorrelation length parameters used to compute cloud cover (yellow marks) are those of the 10 best simulations from the Hybrid experiment; in c-l, 3D-RT (green) and 1D-RT (orange) are both configured with the same 10 best Hybrid configurations.

496 The very best configuration in the Hybrid experiment corresponds to large  $FSD$   
 497 and small  $\ell$  values, close to the best 1D-RT simulations. Accordingly, in the upper panel  
 498 of Figure 4, this best Hybrid configuration corresponds to the yellow marks that are clos-  
 499 est to the 3D LES cloud cover. Comparing fluxes in Figure 4 with those of Figure 3, we  
 500 can see that 1D-RT fluxes, which were previously extremely accurate, are now less so:  
 501 most of the ten best configurations from the Hybrid experiment lead to smaller cloud  
 502 cover than in the 3D LES field, which negatively impacts 1D-RT performances. As for  
 503 3D-RT, it seems to perform better or equally at  $SZA=0^\circ$ , and to be slightly worse than  
 504 before at large solar zenith angles. If we look at the “large  $\ell$ , small  $FSD$ ” configurations  
 505 of the Hybrid experiment, we note that, for a given  $\ell$  value,  $FSD$  is smaller than in the  
 506 pure 3D-RT experiment (compare for instance, yellow square with blue diamond). This  
 507 would lead to more reflective clouds than before, if nothing else changed. However, these  
 508 configurations are also associated with larger  $C_s$  values, which means less intense 3D ef-  
 509 fects, and therefore less reflective clouds at large  $SZA$ . In the end, reflectivity at large  
 510  $SZA$  is largely unchanged with respect to the best 3D-RT of Section 3.2. The fact that  
 511 smaller  $C_s$  values would not be chosen for these “large  $\ell$ , small  $FSD$ ” configurations,  
 512 despite the fact that it would systematically improve fluxes at  $SZA=77^\circ$ , seems to in-  
 513 dicate that fluxes at  $SZA=0^\circ$  would be too large (more intense 3D effects would decrease  
 514 reflectivity at  $SZA=0^\circ$ ) and would become the limiting metrics.

515 Finally, if our hypothesis that a good 3D-RT model should be able to simulate 1D  
 516 fluxes once 3D effects are removed is valid, then the fact that scores are smaller in the  
 517 best 3D-RT configurations of Section 3.2 is a sign of over-fitting. As a consequence, the  
 518 scores of the best Hybrid configurations, and not those of the best 3D-RT, should be used  
 519 to derive estimates of structural errors of the 3D-RT solver (errors that remain once para-  
 520 metric errors are removed). This assumption and its implications will be further discussed  
 521 in Section 5. In the meantime, let us conclude that two 3D-RT configurations are po-  
 522 tential candidate for “best reflectivity”: the “best 3D-RT” configuration and the “best  
 523 Hybrid” one, even if the first one might be there for the wrong reasons.

## 524 4 Compensations between radiative transfer and cloud models

525 In Section 3, we went into details of internal compensations of cloud-geometry as-  
 526 pects within the RT schemes, in a perfect-cloud setup. In the following, we aim at char-  
 527 acterizing error compensations between RT schemes and boundary-layer cloud param-  
 528 eterizations within a SCM. We wonder if errors in convection and cloud schemes may  
 529 arise from the need to compensate cloud-reflectivity errors when RT is inaccurate. In-  
 530 deed, the column (or scene) reflectivity is a combination of physical cloud properties in  
 531 the column (namely, cloud fraction and liquid water content profiles) and cloud reflect-  
 532 ivity (itself a function of cloud geometry). Building on the detailed analysis of Section 3,  
 533 we investigate how cloud-reflectivity errors might propagate through tuning to errors on  
 534 cloud fractions simulated by the SCM to result in accurate scene reflectivity. We thus  
 535 study cloud properties of SCM PPEs produced under radiative constraints, in experi-  
 536 ments that use different RT configurations with varying levels of errors on cloud reflect-  
 537 ivity. Since they yield relatively small cloud-reflectivity errors, experiments based on  
 538 3D-RT configurations are regarded as control experiments, and error compensations are  
 539 examined in equivalent 1D-RT-based experiments, that is, using 1D-RT configurations  
 540 that are the same as the control 3D-RT but without 3D effects.

### 541 4.1 Experimental design

542 The experiments of this section are designed to mimic the tuning protocol often  
 543 followed in GCMs when targeting top-of-the-atmosphere (TOA) radiative fluxes from  
 544 satellite observations. In our idealized SCM version of the protocol, and contrary to Sec-  
 545 tion 3, experiments based on either 1D- or 3D-RT target the same 3D MC fluxes (“true”

radiative fluxes). Also contrary to Section 3, in which the metrics were instantaneous  
 fluxes, here the metrics are hourly-mean fluxes. References are obtained by averaging  
 MC fluxes computed every 5 min on 3D LES outputs. In the SCM, fluxes are calculated  
 on hourly-averaged cloud profiles.

Thirteen parameters that control boundary-layer and cloud parameterizations in  
 LMDZ are varied as in Hourdin et al. (2021) and Hourdin et al. (2023) (see details in  
 Table S1 of Supporting Information). Because we specifically want to investigate compensating  
 errors between radiation and clouds, radiative parameters are not varied in  
 these experiments. This is not necessarily the strategy that would be employed in the  
 tuning of a GCM. *[rev3 2.9 Could you include the radiative transfer parameters in this  
 experiment as well? Could it be a way to mitigate the compensating errors? Naj : j'ai  
 précisé]*

Each tuning experiment consists in 30 iterations. At each iteration, 130 free-parameter  
 vectors are sampled in the NROY space and ARMCU and RICO are simulated using these  
 130 configurations of LMDZ. One RT configuration is then run offline for the two cho-  
 sen hours of the two cumulus cases and each of the 130 SCM configurations, to compute  
 solar reflected fluxes at three solar zenith angles each; in total, 1560 ecRad runs per it-  
 eration. Then one emulator is built for each of the 12 metrics, using their 130 evalua-  
 tions as a learning database. Finally, implausible SCM parameter vectors are ruled out  
 using a threshold that varies from 3 in the first 5 iterations to 2.5 in iterations 6 to 10,  
 to 2 in iterations 11 to 30. The NROY space is thus efficiently narrowed down.

Tolerance to error is set to  $T_{3D} = 3.5 \text{ W}\cdot\text{m}^{-2}$  for all metrics of 3D-RT experiments,  
 and  $T_{1D} = 9 \text{ W}\cdot\text{m}^{-2}$  for all metrics of 1D-RT experiments.  $T_{3D}$  was found by trial and  
 error, seeking the smallest value still yielding non-empty NROY space without falling  
 under  $3.3 \text{ W}\cdot\text{m}^{-2}$ , which is the minimum estimated structural error of 3D-RT (see Sup-  
 porting Information Text S1).  $T_{1D}$  was obtained following another, cheaper strategy: a  
 preliminary experiment was run using an arbitrarily low tolerance to error, and the best  
 simulations were selected amongst the 5 waves that were performed before the NROY  
 space was emptied. Setting  $T_{1D}$  to the errors of these best simulations ensures that it  
 exists configurations satisfying this requirement, but it might be a large overestimate of  
 structural error if only bad configurations were explored during the preliminary exper-  
 iment.

## 4.2 Cloud reflectivity vs. cloud fraction error compensations for too-compact clouds

*[MCD : pas sur de ma traduction Spartacus et Tripleclouds ici. A reprendre a tete reposée] [[j'ai repris]*

In these first experiments, the RT schemes are configured using the “best-3D” con-  
 figuration of Section 3.2 (configuration 8-063 in Figure 2). The 3D-RT simulations use  
 the full configuration while the 1D-RT ones only use overlap and inhomogeneity values.  
 The version of 3D-RT that yields the best reflectivity is regarded as a control, and cloud-  
 reflectivity errors are introduced by using 1D-RT instead of 3D.

For each experiment, maximum and average absolute errors (metrics-wise) of the  
 best-score simulation are retained as a measure of the accuracy of the model. They are  
 presented in Table 2. Note that 3D-RT fluxes are around as accurate as when 3D-RT  
 was run on reference cloud profiles, since tolerance to error in Section 3 was of  $3.5 \text{ W}\cdot\text{m}^{-2}$   
 and the best 3D-RT score was close to 1. On the other hand, flux errors in the 1D-RT  
 experiment are about four times larger than those of 3D-RT. This means that even with  
 the possibility of compensating radiation error with biased cloud properties, the 1D-RT  
 scheme could not produce fluxes as accurate as 3D-RT in these experiments. Note that  
 1D-RT scores are much larger than in perfect cloud experiments but this is partly be-

**Table 2.** Fluxes and cloud fraction errors for the experiments of Section 4.

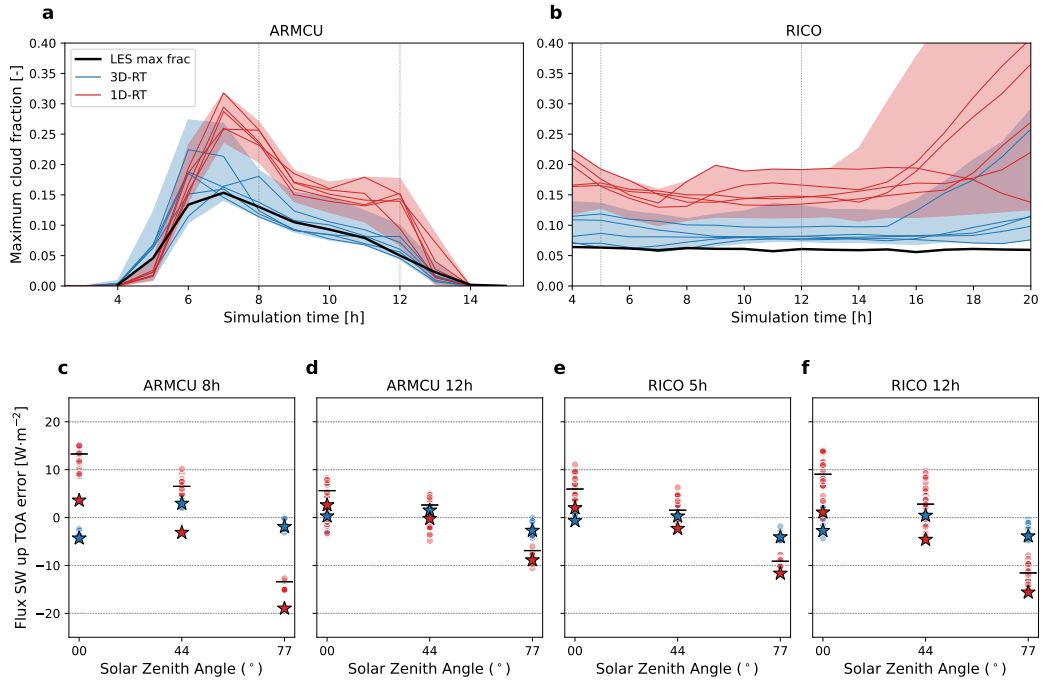
Configuration Best / Solver / Metrics	Shown in Fig.	Flux MxAE ( $\text{W}\cdot\text{m}^{-2}$ )	Flux MAE ( $\text{W}\cdot\text{m}^{-2}$ )	CF MxAE (-)	CF MAE (-)
3D-RT / 3D-RT / 3 SZA	5	3.6480	2.4931	0.0234	0.0156
3D-RT / 1D-RT / 3 SZA	5	14.2253	8.2945	0.1261	0.0941
Hybrid / 3D-RT / 3 SZA	6	3.7890	2.8047	0.0568	0.0444
Hybrid / 1D-RT / 3 SZA	6	14.2117	8.3226	0.1213	0.0841
Hybrid / 3D-RT / $\langle$ SZA $\rangle$	7	0.4360	0.2976	0.0462	0.0368
Hybrid / 1D-RT / $\langle$ SZA $\rangle$	7	1.2210	0.7972	0.0644	0.0482
Hybrid / 1D-RT / 2 SZA	S2	12.6082	8.7031	0.1230	0.0852
1D-RT / 1D-RT / 3 SZA	S2	11.2147	7.4093	0.1085	0.0717

MxAE: Maximum Absolute Error. MAE: Mean Absolute Error. CF: Maximum Cloud Fraction. Mean and max flux error is over the 4 scenes  $\times$  number of SZA. Mean and max CF error is over the 4 scenes.

596 cause in previous experiments it was compared to 1D MC fluxes, whereas it is now com-  
597 pared to 3D MC fluxes.

598 Features of the 30 best simulations of each tuning experiment are analyzed in Fig-  
599 ure 5. The top panel presents temporal evolutions of layer-wise maximum cloud fraction  
600 for ARMCU and RICO test cases. The bottom panel shows errors on the TOA upwelling  
601 fluxes as a function of solar zenith angles, for the four cloud scenes where radiative con-  
602 straints are applied (8th and 12th hour of ARMCU, 5th and 12th hour of RICO). Er-  
603 rors due to neglecting 3D effects in the MC calculations, and errors made by 3D-RT and  
604 1D-RT run on LES mean cloud profiles are also represented in the figures. It shows that  
605 although SCM simulations were constrained to be as close as possible to 3D MC fluxes,  
606 those based on 1D-RT they still closely follow the 1D MC fluxes. This is consistent with  
607 results shown in Figure S2 of Supporting Information [MCD : en dire un peu plus ? Pour  
608 Fredho : C'est la figure trpliclouds qui vise MC-3D en mode perfect clouds sur laque-  
609 lle on voit que tripleclouds ne sait pas compenser les effets 3D et les meilleurs simus ont  
610 des paramètres  $\ell$  et  $FSD$  dans la même région que tripleclouds qui vise du 1D. La seule  
611 différence est la tolérance à l'erreur de l'expérience qui doit être augmenté jusqu'à une  
612 dizaine de  $\text{W}\cdot\text{m}^{-2}$  soit l'ordre de grandeur des effets 3D] 3D-RT fluxes however match  
613 the 3D-MC reference for almost every angle, including those that were not constrained  
614 (not shown). In the upper plots, we see that, in addition to better simulating solar ra-  
615 diation, the 30 best SCM simulations of the 3D-RT experiment exhibit cloud fraction  
616 evolutions that are compatible with reference values for ARMCU, and close to the LES  
617 ones for RICO. Conversely, the 30 best 1D-RT simulations persistently overestimate max-  
618 imum cloud fraction for both cases.

619 We interpret this as a sign of compensating errors due to wrong cloud reflectivity  
620 in the 1D-RT scheme. Indeed, plots in the lower row of Figure 5 show that the largest  
621 errors in 1D-RT simulations run on LES mean cloud profiles (red stars) are systemat-  
622 ically the underestimation of reflected solar radiation at  $\text{SZA}=77^\circ$ . The absolute error  
623 is the largest for the 8th hour of ARMCU where it reaches  $-20.5 \text{ W}\cdot\text{m}^{-2}$ . The tuning pro-  
624 tocol, which rejects configurations when the metric-wise maximum error is too large, sys-  
625 tematically selects configurations where the error at  $\text{SZA}=77^\circ$  is smaller than that of “per-  
626 fect clouds”, thanks to an increase of cloud fraction that compensates for a lack of cloud  
627 reflectivity. This increase of cloud fraction is limited by the fact that it also increases  
628 reflected radiation at  $\text{SZA}=0^\circ$ , which cannot be too large either. The tuning finally se-  
629 lects configurations for which the bias shows the same absolute value for the maximum  
630 underestimation (about  $-15 \text{ W}\cdot\text{m}^{-2}$  at  $\text{SZA}=77^\circ$  for the 8th hour of ARMCU) and max-



**Figure 5.** Clouds and radiation for the 30 best simulations for two experiments calibrating LMDZ parameters under radiative constraints using either **Tripleclouds** (red) or **Spartacus** (blue). Both solvers use parameters of the best **Spartacus** configuration from Section 3.2 (**Spartacus** 8-063). Left: ARMCU. Right: RICO. a-b: maximum cloud fraction (shadings represent the 30 best envelope, lines show the 5 best). c-f: error in upwelling TOA fluxes as a function of solar zenith angle, compared to MC 3D, for the four cloud scenes: ARMCU 8h (c), ARMCU 12h (d), RICO 5h (e), RICO 12h (f).

631 inum overestimation (about  $15 \text{ W}\cdot\text{m}^{-2}$  at  $\text{SZA}=0^\circ$  for the 8th hour of ARMCU and 12th  
632 hour of RICO).

633 In the control experiment based on 3D-RT, angular dependence of reflected radi-  
634 ation is much better represented, avoiding compromises between errors at high and low  
635 zenith angles. Overall, there is no sign of cloud–radiation compensating errors when the  
636 3D-RT scheme is used: The best SCM configurations obtained with a 3D-RT-based tun-  
637 ing produce at the same time (i) fluxes that is as good on average as in the perfect-cloud  
638 experiments, and (ii) maximum cloud fractions that correspond reasonably well to those  
639 of the LES.

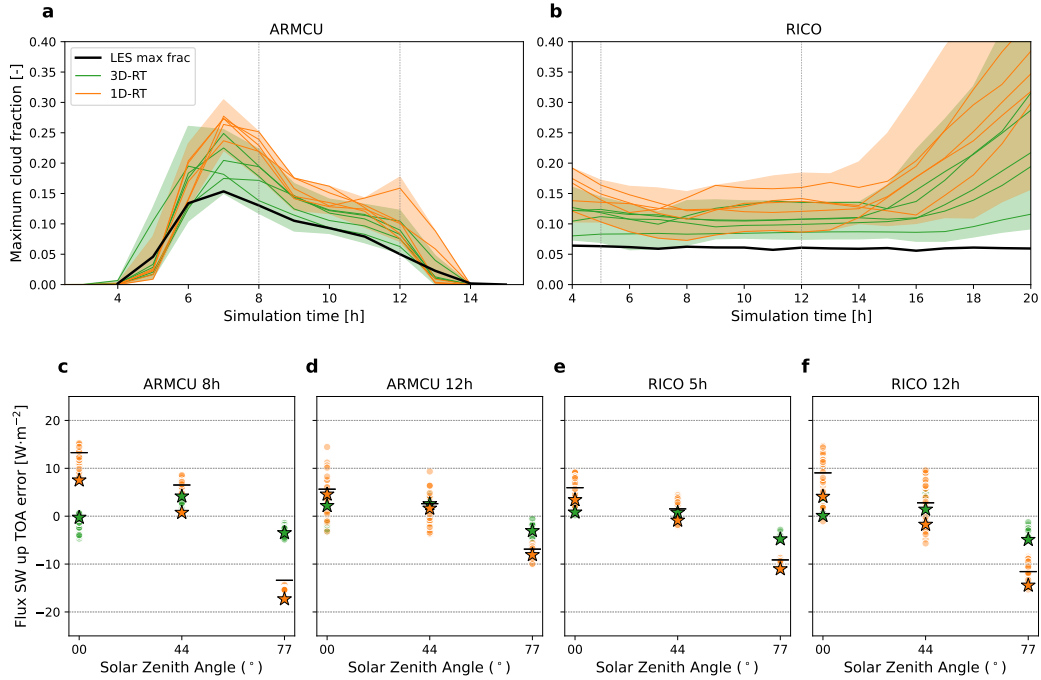
640 Here, the only difference between the two experiments is the removal of 3D effects  
641 in the 1D-RT-based experiment. The intuitive conclusion at this stage would be that,  
642 in the 1D-RT-based experiment, cloud fractions simulated by the SCM are overestimated  
643 to compensate for lack of 3D effects. However, we saw in Section 3 that removing 3D  
644 effects from the best 3D-RT configurations led to bad 1D-RT configurations, that is, to  
645  $(\ell, FSD)$  values which result in an underestimation of cloud reflectivity. This means that  
646 in this 1D-RT-based experiment, cloud reflectivity is not only underestimated at large  
647 SZA because 3D effects are neglected, but also because the other geometry parameters  
648 are not adequately chosen. Even if the difference between the control and 1D-RT exper-  
649 iments are indeed the strict removal of 3D effects, it would be inaccurate to conclude that  
650 it is the lack of 3D effects that leads to overestimating cloud fraction, when it might in-  
651 stead be due, to some extent, to parametric errors.

### 652 **4.3 Cloud reflectivity vs. cloud fraction error compensations due to lack** 653 **of 3D effects**

654 In order to quantify the impact of neglecting 3D effects in a more robust way, ad-  
655 ditional tuning experiments are performed. First, the “best Hybrid” configuration is used  
656 instead of the “best 3D-RT” one. *[ rev3 2.10 Reading this, it seems that the “best Triple-  
657 clouds” was not used in sections 4.2. Naj : ben oui c’est le cas. Ca veut dire qu’il a cru  
658 qu’on avait pris la best TC au coup d’avant... j’espère que c’est plus clair dans la nou-  
659 velle version]* The rationale is that in the perfect-cloud setup of Section 3, the overlap  
660 parameter of the “best Hybrid” configuration yields cloud cover estimates that are close  
661 to the 3D LES, which in turn leads to more accurate 1D-RT reflectivity. Therefore, us-  
662 ing the “best Hybrid” configuration, we expect the remaining 1D-RT cloud-reflectivity  
663 errors to be smaller than with the “best 3D-RT” and hence, if our hypothesis is correct,  
664 cloud fraction overestimation to be less pronounced. *[rev3, remarques generales sur ce  
665 paragraphe : “I read the first part of this section multiple times without completely un-  
666 derstanding the objective. The title of this section is not clear either. ”]*

667 This is verified in Figure 6, which is the same as Figure 5 but with both solvers us-  
668 ing the “best Hybrid” configuration instead of “best 3D-RT”. The score of the 3D-RT  
669 experiment is a little bit larger than before, suggesting that the degradation introduced  
670 in changing the configuration is not compensated by clouds. The time evolution of max-  
671 imum cloud fractions in the 3D-RT experiments resembles those of Figure 5, which con-  
672 firms that changing cloud-geometry parameters to a slightly less good configuration does  
673 not significantly impact 3D-RT fluxes. In the 1D-RT experiment, the radiative score was  
674 not improved by much. Yet, the evolution of maximum cloud fraction is closer to the LES:  
675 To reach the same radiative accuracy with the 1D-RT scheme, that is, the same scene  
676 reflectivity, cloud fractions did not have to be as large as before. Our interpretation is  
677 that the error in cloud reflectivity that is compensated by increasing cloud fraction is  
678 smaller: perfect-cloud 1D-RT errors at  $\text{SZA } 77^\circ$  are smaller with this best Hybrid con-  
679 figuration, around  $-18.8 \text{ W}\cdot\text{m}^{-2}$  (against  $-20.5$ ) at the 8th hour of ARMCU.

680 Yet, maximum cloud fraction is still overestimated in the 1D-RT-based experiment,  
681 compared to 3D-RT. This might be for two reasons: in the best Hybrid configuration,

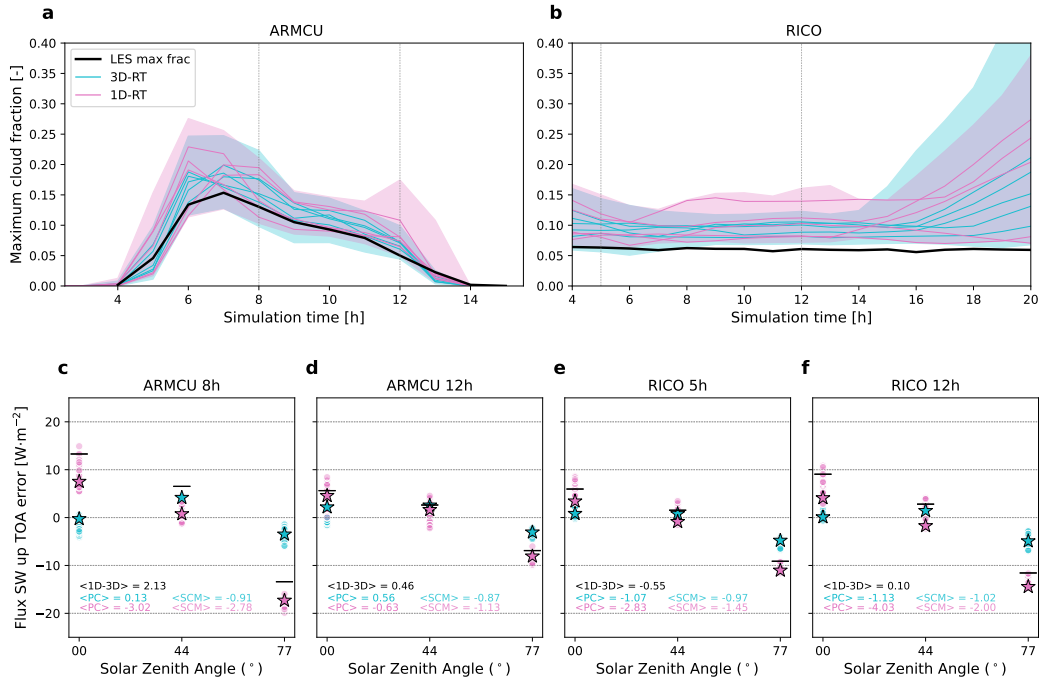


**Figure 6.** Clouds and radiation for the 30 best simulations for two experiments calibrating LMDZ parameters under radiative constraints using either `Tripleclouds` (orange) or `Spartacus` (green). Same as Figure 5 except that both solvers now use parameters of the best Hybrid configuration from Section 3.3 (Hybrid 16-092).

682 overlap parameter is similar to best 1D-RT values, but the associated FSD is larger (for  
 683 for the same  $\ell$ ) which means that clouds are less reflective; again, in 3D-RT, this overlap-  
 684 heterogeneity combination is balanced by 3D effects, but in 1D-RT a large  $FSD$  leads  
 685 to an underestimated cloud reflectivity. Hence, maximum cloud fraction might again be  
 686 overestimated in compensation of this parametric error. The second reason is that, in-  
 687 dependently of heterogeneity, accurate 3D fluxes cannot be obtained simultaneously for  
 688 all solar angles, for the right clouds, if 3D effects are not represented in the model.

689 A final pair of experiments is done in which we change the target metrics to an average  
 690 over eight solar zenith angles, from  $0^\circ$  to  $77^\circ$  with step  $11^\circ$ . As 3D effects change  
 691 sign when sun zenith angle increases, they partially cancel out on average. If lack of 3D  
 692 effects is indeed the error that is compensated by cloud fractions in the 1D-RT exper-  
 693 iment, then it should be less pronounced when targeting an angular average of fluxes,  
 694 in which the expression of 3D effects is almost null. It is the unweighted arithmetic av-  
 695 erage that is considered here, as the simplest way to test our assumption, although more  
 696 complex averages that better represent solar angle distributions on Earth might be more  
 697 relevant in a real GCM tuning exercise.

698 Results are presented in Figure 7. With the average flux metric, cloud fraction evolu-  
 699 tion is the same in the two experiments, using 3D-RT or 1D-RT. It means that 1D-  
 700 RT errors ( $FSD$  too large or lack of 3D effects) are no longer compensated by cloud frac-  
 701 tions. Looking at the detailed flux errors, we see that 1D-RT flux at  $77^\circ$  are more neg-  
 702 atively biased compared to previous experiments, in agreement with the fact that this  
 703 particular value is no longer constrained. It confirms that maximum cloud fractions were  
 704 previously overestimated to compensate for lack of 3D effects in particular at this large



**Figure 7.** Clouds and radiation for the 30 best simulations for two experiments calibrating LMDZ parameters under radiative constraints using either *Tripleclouds* (pink) or *Spartacus* (cyan). Same as Figure 6 except that the radiative constraint is now the flux averaged over solar zenith angles, instead of flux values taken at three different SZA. c-f: In addition to flux errors as a function of SZA (markers), SZA-mean flux errors are written:  $\langle 1D-3D \rangle$  is the average 3D effects,  $\langle PC \rangle$  is for Perfect Cloud experiments (associated with star markers),  $\langle SCM \rangle$  is for SCM-tuned experiments (associated with bullet markers).

705 solar zenith angle; because increasing cloud fraction also increases reflected flux at all  
 706 other angles, in 1D-RT experiments, flux errors could not be entirely compensated by  
 707 clouds: increasing cloud fraction too much would have led to even larger flux overesti-  
 708 mation at  $\text{SZA}=0^\circ$ , which was not acceptable given the metrics tolerance to error. To  
 709 support this statement, a complementary experiment where the SCM was tuned based  
 710 on 1D-RT but using only constraints at  $\text{SZA}=44^\circ$  and  $\text{SZA}=77^\circ$  is presented in Support-  
 711 ing Information (Figure S2).

712 This last experiment, targeting average flux, is an important result for climate mod-  
 713 eling: it confirms that 3D radiative effects could, for a large part, self compensate when  
 714 averaged over diurnal cycle, seasons or latitudes. It also proposes a simple compromise  
 715 to tune a model with a 1D RT code, against LES/MC reference simulations. However,  
 716 these results also suggest the possibility of compensating errors at high latitudes, where  
 717 cloud fraction might be overestimated to compensate for lack of 3D radiative effects, if  
 718 the model was tuned targeting high-latitude metrics.

## 719 5 Conclusion and discussion

720 By tuning parameterizations of a climate model in SCM mode targeting TOA up-  
 721 ward solar fluxes, we evidenced mechanisms of error ~~compensations~~ compensation both  
 722 inside the RT model itself and in interaction with the cloud physics simulations.

723 Considering the RT computation itself, two kinds of ~~compensations~~ compensation  
 724 were identified. ~~First we show~~ We showed first how the cloud vertical overlap and in-cloud  
 725 heterogeneities can compensate each other inside the RT solver to simulate fairly accu-  
 726 rate 1D fluxes, a compensation already known and studied in 1D-RT as well as 3D-RT  
 727 (Villefranque et al., 2021). Secondly, when tuning 3D-RT computations targeting 3D-  
 728 MC simulations, we showed how the best performing simulations selected favour param-  
 729 eters which underestimate the 1D reflectivity, in the sense that 1D-RT simulations run  
 730 with the same parameter values underestimate the reflected solar radiation at TOA. This  
 731 underestimation is consistent with the fact that the total cloud cover diagnosed in Ecrad  
 732 by the exponential-random scheme is underestimated compared to the LES cloud cover.  
 733 One way to interpret these results would be to say that errors in the computation of 3D  
 734 effects in the 3D-RT *Spartacus* solver are compensated by an underestimation of cloud  
 735 cover or 1D reflectivity.

736 This way of presenting the second mechanism for error compensation relies how-  
 737 ever on the hypothesis that 3D effects would be added on the results of a 1D-RT com-  
 738 putation. In this case, it would be legitimate to privilege values that are optimal both  
 739 in 1D-RT and 3D-RT mode, in order to fully avoid compensating errors between 3D effects  
 740 and cloud cover or reflectivity. Actually, when running the 3D-RT solver, there is no real  
 741 distinction between non-3D and 3D effects. It is the full radiation which is computed and  
 742 affected by the parameterization of 3D effects. It may ~~but so not~~ be a surprise ~~so~~ that  
 743 the tuning done in 1D-RT targeting 1D-MC, considering the number of numerous ap-  
 744 proximations and unavoidable error compensation inside the 1D-RT scheme itself, priv-  
 745 ilege parameter values which are not optimal when targeting 3D-MC fluxes using a more  
 746 physical radiative transfer model. From this perspective, we consider free parameters as  
 747 *effective parameters* that are meaningful only within their respective models and this con-  
 748 straint relaxes. Related to this, the cloud cover that is directly determined by the decor-  
 749 relation length free parameter has no reason to be exactly the same in 1D and 3D-RT  
 750 solver, nor between RT SCM computations and MC references. ~~Here also, this diagnostics~~  
 751 ~~only documents approximately how an effective cloud cover is related to the computation~~  
 752 ~~of RT, accounting for the effect of cloud overlap.~~ Thereby, constraining the 3D-RT solver  
 753 to well simulate 1D fluxes when 3D radiative effects are removed is not a straightforward  
 754 choice. ~~On the one hand, there~~ There is no clear or definite reason to privilege param-  
 755 eter values that would be the best both in 1D-RT and 3D-RT computations. ~~On the other~~

756 ~~hand,~~ Note however that the fact that we were able to find parameter values for which  
 757 RT computations match reasonably well, even if not optimally, both the 1D-MC and 3D-  
 758 MC fluxes in 1D-RT computation and 3D-RT computation respectively, together with  
 759 a reasonable diagnostic of the cloud cover compared to LES, is a strong indication that  
 760 the “reasonably good simulations” were obtained for the “reasonably good physics”. Be-  
 761 yond these specific results, this study illustrates the step forward allowed by being able  
 762 to start address error compensation as a scientific question, ~~but opens as well~~ opening  
 763 in the mean time a number of ~~open~~ questions that modelers should have in mind when  
 764 using, evaluating and tuning their models.

765 We show in a second part that structural errors in the 1D-RT model can be com-  
 766 pensated by errors in cloud properties when TOA radiative fluxes are targeted in a tun-  
 767 ing process. ~~Maximum~~ The maximum cloud fractions obtained after tuning of the free  
 768 parameters of the cloud and boundary layer model targeting the 3D-MC flux ~~computations,~~  
 769 ~~are~~ systematically overestimated to compensate two possible errors in the 1D-RT. The  
 770 first one concerns the choice of cloud geometry parameter values, that lead to too com-  
 771 pact and heterogeneous cloud from 1D-RT perspective when selected to perform well in  
 772 the 3D-RT solver. The second one is a compensation of the structural error ~~of~~ due to  
 773 not accounting for 3D radiative effects, that lead to a strong underestimation of cloud  
 774 reflectivity at high SZA. This result provides a novel argument in favor of modelling 3D  
 775 radiative effects in climate models: even if they were small on average and had a weak  
 776 feedback on circulations and climate, we have shown that systematic errors in radiative  
 777 transfer can generate systematic errors in other components of the model through tun-  
 778 ing. A better radiative transfer model might remove the risk of degrading the clouds rep-  
 779 resentations to obtain ~~the~~ right fluxes, in particular at high latitudes where the occur-  
 780 rence of high zenith angles is larger.

781 This demonstration is made in an idealized configuration, and our results should  
 782 not be directly extrapolated to 3D coupled climate models. Indeed, in the SCM setup  
 783 considered here, only shallow convection and cloud parameterizations can compensate  
 784 structural radiative errors, whereas much more processes are at work in a 3D GCM, which  
 785 can result in other compensating errors. ~~Also, radiative fluxes were the only constraint,~~  
 786 ~~but~~ Couvreux et al. (2021); Hourdin et al. (2021) Couvreux et al. (2021) and Hourdin et al. (2021)  
 787 claim that compensating errors can be ~~prevented or at least~~ limited by process-based tun-  
 788 ing in SCM mode before tuning the full GCM. With this strategy, constraints can be set  
 789 directly on cloud properties to rule out model configurations that yield wrong cloud frac-  
 790 tions. Note however that tuning towards radiative targets while preventing clouds from  
 791 compensating radiation errors might generate compensating errors elsewhere in the sys-  
 792 tem.

793 The main structural error under study here is known and introduced intentionally,  
 794 using the 1D-RT instead of 3D-RT configuration, that is generally not the case when eval-  
 795 uating and tuning a climate model. This allowed us to demonstrate the potential of the  
 796 History Matching tuning method for the quantitative investigation of compensating er-  
 797 rors at work in climate models. Through tuning we explore model configurations and  
 798 resulting simulated climates in a parameter space reduced under a set of chosen constraints.  
 799 This allows us to disentangle parametric from structural errors. Notably, when no set  
 800 of parameters can be found for which all simulated metrics comply with user require-  
 801 ments, it indicates that structural errors are larger than tolerated errors, and hence that  
 802 the model is incomplete. This is a powerful way to detect where the development effort  
 803 have to be made to improve the simulation of the key climate characteristics by the model.  
 804 On the other side, when simulated metrics do comply with prescribed requirements, re-  
 805 sulting perturbed parameter ensembles of simulations (PPE) can be used to investigate  
 806 compensating errors, better understand the model and its physics through global sen-  
 807 sitivity studies, and quantify parametric uncertainties on various aspects of climate.

808 Finally, the tuning tool used here, High-Tune:Explorer, is based on machine learn-  
 809 ing techniques: predictive Gaussian Processes are trained on a small amount of simu-  
 810 lated data and are then able to emulate the model’s response much faster than the ac-  
 811 tual model. Thanks to this approach, the model’s high-dimensional parameter space can  
 812 be explored and shrunken efficiently. Machine learning is here at the service of physics;  
 813 it helps saving computing time but not at the expense of the physical consistency of the  
 814 model. This consistency is crucial for our confidence in climate projections and to keep  
 815 using models as tools to better understand climate. In the same spirit, we believe that  
 816 research that aims at better understanding climate models and the act of modeling it-  
 817 self is a crucial aspect of climate sciences.

## 818 Open Research Section

819 High-Tune Explorer (htexplo) and LMDZ are available through the open source  
 820 version control system “subversion” (svn). htexplo is distributed under the GPL-v3 li-  
 821 cense, and LMDZ is distributed under the CeCILL version 2 license. The htexplo release  
 822 used in the study can be downloaded through `svn checkout http://svn.lmd.jussieu.fr/HighTune`  
 823 `-r 568`. The LMDZ release used in the study can be configured and installed directly  
 824 on Linux machines with an installation bash script `https://lmdz.lmd.jussieu.fr/pub/install_lmdz.sh`  
 825 run with `as bash install_lmdz.sh -SCM -v 20250337.trunk` The ecRad offline pack-  
 826 age is freely available under the terms of the Apache License Version 2.0. The release  
 827 used in this study corresponds to commit fa642e, which is based on version v1.6-beta.  
 828 A tar file of the htexplo, LMDZ and ecRad codes used as well as the data that supports  
 829 this research, the results of the SCM simulations, as well as the scripts for visualization  
 830 WILL BE MADE AVAILABLE ON A DOI IF THE PAPER IS ACCEPTED FOR PUB-  
 831 LICATION. The corresponding DOIs will be provided during galley proofs by placeholder  
 832 “IPSL data catalog.”

## 833 Acknowledgments

834 The PhD scholarship of the first author was funded by Institut Pierre-Simon Laplace.  
 835 This work was supported by the CNRS and the GDR DEPHY. The authors wish to thank  
 836 the LMDZ team and the htexplo team for their investment in the development and main-  
 837 tainance of community tools.

## 838 References

- 839 Barker, H. W., Stephens, G. L., & Fu, Q. (1999). The sensitivity of domain-averaged  
 840 solar fluxes to assumptions about cloud geometry. *Quarterly Journal of the*  
 841 *Royal Meteorological Society*, 125(558), 2127–2152. Retrieved from [https://](https://rmets.onlinelibrary.wiley.com/doi/abs/10.1002/qj.49712555810)  
 842 [rmets.onlinelibrary.wiley.com/doi/abs/10.1002/qj.49712555810](https://rmets.onlinelibrary.wiley.com/doi/abs/10.1002/qj.49712555810) doi:  
 843 10.1002/qj.49712555810
- 844 Brown, A. R., Cederwall, R. T., Chlond, A., Duynkerke, P. G., Golaz, J.-C.,  
 845 Khairoutdinov, M., ... Stevens, B. (2002). Large-eddy simulation of  
 846 the diurnal cycle of shallow cumulus convection over land. *Quarterly*  
 847 *Journal of the Royal Meteorological Society*, 128(582), 1075–1093. doi:  
 848 10.1256/003590002320373210
- 849 Couvreur, F., Hourdin, F., Williamson, D., Roehrig, R., Volodina, V., Villefranque,  
 850 N., ... Xu, W. (2021). Process-Based Climate Model Development Harnessing  
 851 Machine Learning: I. A Calibration Tool for Parameterization Improvement.  
 852 *Journal of Advances in Modeling Earth Systems*, 13(3), e2020MS002217. doi:  
 853 10.1029/2020MS002217
- 854 Dephy. (2020). *Dephy-scm: Single-column model standards and case drivers*. Re-  
 855 trieved from <https://github.com/gdr-dephy/dephy-scm>
- 856 Hogan, R. J., & Bozzo, A. (2018). A Flexible and Efficient Radiation Scheme for

- 857 the ECMWF Model. *Journal of Advances in Modeling Earth Systems*, 10(8),  
 858 1990–2008. doi: 10.1029/2018MS001364
- 859 Hogan, R. J., Fielding, M. D., Barker, H. W., Villefranche, N., & Schäfer, S. A. K.  
 860 (2019). Entrapment: An Important Mechanism to Explain the Shortwave  
 861 3D Radiative Effect of Clouds. *Journal of the Atmospheric Sciences*, 76(7),  
 862 2123–2141. doi: 10.1175/JAS-D-18-0366.1
- 863 Hogan, R. J., & Illingworth, A. J. (2000, October). Deriving cloud overlap statistics  
 864 from radar. *Quarterly Journal of the Royal Meteorological Society*, 126(569),  
 865 2903–2909. doi: 10.1002/qj.49712656914
- 866 Hogan, R. J., Schäfer, S. A. K., Klinger, C., Chiu, J. C., & Mayer, B. (2016).  
 867 Representing 3-D cloud radiation effects in two-stream schemes: 2. Matrix  
 868 formulation and broadband evaluation. *Journal of Geophysical Research:  
 869 Atmospheres*, 121(14), 8583–8599. doi: 10.1002/2016JD024875
- 870 Hogan, R. J., & Shonk, J. K. P. (2013, February). Incorporating the Effects of 3D  
 871 Radiative Transfer in the Presence of Clouds into Two-Stream Multilayer Ra-  
 872 diation Schemes. *Journal of the Atmospheric Sciences*, 70(2), 708–724. doi:  
 873 10.1175/JAS-D-12-041.1
- 874 Hourdin, F., Ferster, B., Deshayes, J., Mignot, J., Musat, I., & Williamson, D.  
 875 (2023, July). Toward machine-assisted tuning avoiding the underestimation of  
 876 uncertainty in climate change projections. *Science Advances*, 9(29), eadf2758.  
 877 doi: 10.1126/sciadv.adf2758
- 878 Hourdin, F., Grandpeix, J.-Y., Rio, C., Bony, S., Jam, A., Cheruy, F., ... Roehrig,  
 879 R. (2013, May). LMDZ5B: The atmospheric component of the IPSL climate  
 880 model with revisited parameterizations for clouds and convection. *Climate  
 881 Dynamics*, 40(9-10), 2193–2222. doi: 10.1007/s00382-012-1343-y
- 882 Hourdin, F., Jam, A., Rio, C., Couvreux, F., Sandu, I., Lefebvre, M.-P., ... Idelkadi,  
 883 A. (2019). Unified Parameterization of Convective Boundary Layer Transport  
 884 and Clouds With the Thermal Plume Model. *Journal of Advances in Modeling  
 885 Earth Systems*, 11(9), 2910–2933. doi: 10.1029/2019MS001666
- 886 Hourdin, F., Mauritsen, T., Gettelman, A., Golaz, J.-C., Balaji, V., Duan, Q., ...  
 887 Williamson, D. (2017, March). The Art and Science of Climate Model Tun-  
 888 ing. *Bulletin of the American Meteorological Society*, 98(3), 589–602. doi:  
 889 10.1175/BAMS-D-15-00135.1
- 890 Hourdin, F., Rio, C., Grandpeix, J.-Y., Madeleine, J.-B., Cheruy, F., Rochetin,  
 891 N., ... Ghattas, J. (2020). LMDZ6A: The Atmospheric Component of  
 892 the IPSL Climate Model With Improved and Better Tuned Physics. *Jour-  
 893 nal of Advances in Modeling Earth Systems*, 12(7), e2019MS001892. doi:  
 894 10.1029/2019MS001892
- 895 Hourdin, F., Williamson, D., Rio, C., Couvreux, F., Roehrig, R., Villefranche, N.,  
 896 ... Volodina, V. (2021). Process-Based Climate Model Development Harness-  
 897 ing Machine Learning: II. Model Calibration From Single Column to Global.  
 898 *Journal of Advances in Modeling Earth Systems*, 13(6), e2020MS002225. doi:  
 899 10.1029/2020MS002225
- 900 Jam, A., Hourdin, F., Rio, C., & Couvreux, F. (2013, June). Resolved Versus  
 901 Parametrized Boundary-Layer Plumes. Part III: Derivation of a Statistical  
 902 Scheme for Cumulus Clouds. *Boundary-Layer Meteorology*, 147(3), 421–441.  
 903 doi: 10.1007/s10546-012-9789-3
- 904 Konsta, D., Dufresne, J.-L., Chepfer, H., Vial, J., Koshiro, T., Kawai, H., ... Ogura,  
 905 T. (2022). Low-Level Marine Tropical Clouds in Six CMIP6 Models Are Too  
 906 Few, Too Bright but Also Too Compact and Too Homogeneous. *Geophysical  
 907 Research Letters*, 49(11), e2021GL097593. doi: 10.1029/2021GL097593
- 908 Lac, C., Chaboureau, J.-P., Masson, V., Pinty, J.-P., Tulet, P., Escobar, J., ...  
 909 Wautelet, P. (2018, May). Overview of the Meso-NH model version 5.4 and  
 910 its applications. *Geoscientific Model Development*, 11(5), 1929–1969. doi:  
 911 10.5194/gmd-11-1929-2018

- 912 Lafore, J. P., Stein, J., Asencio, N., Bougeault, P., Ducrocq, V., Duron, J., ... Vilà-  
 913 Guerau de Arellano, J. (1998, January). The Meso-NH Atmospheric Simula-  
 914 tion System. Part I: Adiabatic formulation and control simulations. *Annales*  
 915 *Geophysicae*, *16*(1), 90–109. doi: 10.1007/s00585-997-0090-6
- 916 Madeleine, J.-B., Hourdin, F., Grandpeix, J.-Y., Rio, C., Dufresne, J.-L., Vignon,  
 917 E., ... Bonazzola, M. (2020). Improved Representation of Clouds in the  
 918 Atmospheric Component LMDZ6A of the IPSL-CM6A Earth System Model.  
 919 *Journal of Advances in Modeling Earth Systems*, *12*(10), e2020MS002046. doi:  
 920 10.1029/2020MS002046
- 921 McKee, T. B., & Cox, S. K. (1974). Scattering of Visible Radiation by Finite  
 922 Clouds. *Journal of the Atmospheric Sciences*, *31*(7), 1885–1892. doi:  
 923 10.1175/1520-0469(1974)031<1885:SOVRBF>2.0.CO;2
- 924 Meador, W. E., & Weaver, W. R. (1980, March). Two-Stream Approxima-  
 925 tions to Radiative Transfer in Planetary Atmospheres: A Unified Descrip-  
 926 tion of Existing Methods and a New Improvement. , *37*(3), 630–643. doi:  
 927 10.1175/1520-0469(1980)037<0630:TSATRT>2.0.CO;2
- 928 Nam, C., Bony, S., Dufresne, J.-L., & Chepfer, H. (2012). The ‘too few, too bright’  
 929 tropical low-cloud problem in CMIP5 models. *Geophysical Research Letters*,  
 930 *39*(21). doi: 10.1029/2012GL053421
- 931 Pincus, R., Barker, H. W., & Morcrette, J.-J. (2003, July). A fast, flexible, ap-  
 932 proximate technique for computing radiative transfer in inhomogeneous cloud  
 933 fields. *Journal of Geophysical Research: Atmospheres*, *108*(D13), n/a–n/a. doi:  
 934 10.1029/2002JD003322
- 935 Schäfer, S. A. K., Hogan, R. J., Klinger, C., Chiu, J. C., & Mayer, B. (2016). Rep-  
 936 resenting 3-D cloud radiation effects in two-stream schemes: 1. Longwave con-  
 937 siderations and effective cloud edge length. *Journal of Geophysical Research:*  
 938 *Atmospheres*, *121*(14), 8567–8582. doi: 10.1002/2016JD024876
- 939 Shonk, J. K. P., & Hogan, R. J. (2008, June). Tripleclouds: An Efficient Method  
 940 for Representing Horizontal Cloud Inhomogeneity in 1D Radiation Schemes by  
 941 Using Three Regions at Each Height. *Journal of Climate*, *21*(11), 2352–2370.  
 942 doi: 10.1175/2007JCLI1940.1
- 943 Shonk, J. K. P., Hogan, R. J., Edwards, J. M., & Mace, G. G. (2010, July). Effect  
 944 of improving representation of horizontal and vertical cloud structure on the  
 945 Earth’s global radiation budget. Part I: review and parametrization. *Quarterly*  
 946 *Journal of the Royal Meteorological Society*, n/a–n/a. doi: 10.1002/qj.647
- 947 vanZanten, M. C., Stevens, B., Nuijens, L., Siebesma, A. P., Ackerman, A. S., Bur-  
 948 net, F., ... Wyszogrodzki, A. (2011). Controls on precipitation and cloudiness  
 949 in simulations of trade-wind cumulus as observed during RICO. *Journal of*  
 950 *Advances in Modeling Earth Systems*, *3*(2). doi: 10.1029/2011MS000056
- 951 Várnai, T., & Davies, R. (1999). Effects of cloud heterogeneities on shortwave radi-  
 952 ation: Comparison of cloud-top variability and internal heterogeneity. *Journal*  
 953 *of the Atmospheric Sciences*, *56*(24), 4206–4224.
- 954 Vernon, I., Goldstein, M., & Bower, R. G. (2010, December). Galaxy Formation: A  
 955 Bayesian Uncertainty Analysis. *Bayesian Analysis*, *05*(04). doi: 10.1214/10  
 956 -ba524
- 957 Villefranque, N., Blanco, S., Couvreux, F., Fournier, R., Gautrais, J., Hogan, R. J.,  
 958 ... Williamson, D. (2021). Process-Based Climate Model Development Har-  
 959 nassing Machine Learning: III. The Representation of Cumulus Geometry and  
 960 Their 3D Radiative Effects. *Journal of Advances in Modeling Earth Systems*,  
 961 *13*(4), e2020MS002423. doi: 10.1029/2020MS002423
- 962 Villefranque, N., Fournier, R., Couvreux, F., Blanco, S., Cornet, C., Eymet, V., ...  
 963 Tregan, J.-M. (2019). A Path-Tracing Monte Carlo Library for 3-D Radiative  
 964 Transfer in Highly Resolved Cloudy Atmospheres. *Journal of Advances in*  
 965 *Modeling Earth Systems*, *11*(8), 2449–2473. doi: 10.1029/2018MS001602
- 966 Webb, M., Senior, C., Bony, S., & Morcrette, J. J. (2001, September). Combin-

- 967 ing ERBE and ISCCP data to assess clouds in the Hadley Centre, ECMWF  
968 and LMD atmospheric climate models. *Climate Dynamics*, 17(12), 905–922.  
969 (WOS:000171263400001) doi: 10.1007/s003820100157
- 970 Williamson, D., Goldstein, M., Allison, L., Blaker, A., Challenor, P., Jackson, L.,  
971 & Yamazaki, K. (2013, October). History matching for exploring and re-  
972 ducing climate model parameter space using observations and a large per-  
973 turbed physics ensemble. *Climate Dynamics*, 41(7-8), 1703–1729. doi:  
974 10.1007/s00382-013-1896-4
- 975 Yamada, T. (1983, January). Simulations of Nocturnal Drainage Flows by a q2l Tur-  
976 bulence Closure Model. *Journal of the Atmospheric Sciences*, 40(1), 91–106.  
977 doi: 10.1175/1520-0469(1983)040<0091:SONDFB>2.0.CO;2