



HAL
open science

Scientific and human errors in a snow model intercomparison

Cecile Menard, Richard Essery, Gerhard Krinner, Gabriele Arduini, Paul Bartlett, Aaron Boone, Claire Brutel-Vuilmet, Eleanor Burke, Matthias Cuntz, Yongjiu Dai, et al.

► **To cite this version:**

Cecile Menard, Richard Essery, Gerhard Krinner, Gabriele Arduini, Paul Bartlett, et al.. Scientific and human errors in a snow model intercomparison. Bulletin of the American Meteorological Society, 2021, 102 (1), pp.E61-E79. 10.1175/BAMS-D-19-0329.1 . hal-03013971v1

HAL Id: hal-03013971

<https://hal.science/hal-03013971v1>

Submitted on 19 Nov 2020 (v1), last revised 7 Jan 2022 (v2)

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Scientific and human errors in a snow model intercomparison



Cecile B. Menard¹, Richard Essery¹, Gerhard Krinner², Gabriele Arduini³, Paul Bartlett⁴, Aaron Boone⁵, Claire Brutel-Vuilmet², Eleanor Burke⁶, Matthias Cuntz⁷, Yongjiu Dai⁸, Bertrand Decharme⁴, Emanuel Dutra⁹, Xing Fang¹⁰, Charles Fierz¹¹, Yeugeniy Gusev¹², Stefan Hagemann¹³, Vanessa Haverd¹⁴, Hyungjun Kim¹⁵, Matthieu Lafaysse¹⁶, Thomas Marke¹⁷, Olga Nasonova¹², Tomoko Nitta¹⁵, Masashi Niwano¹⁸, John Pomeroy¹⁰, Gerd Schädler¹⁹, Vladimir Semenov²⁰, Tatiana Smirnova²¹, Ulrich Strasser¹⁷, Sean Swenson²², Dmitry Turkov²³, Nander Wever^{24,11}, Hua Yuan⁸

¹School of Geosciences, University of Edinburgh, Edinburgh, United Kingdom

²CNRS, Université Grenoble Alpes, Institut de Géosciences de l'Environnement (IGE), Grenoble, France

³European Centre for Medium-Range Weather Forecasts (ECMWF), Reading, United Kingdom

⁴Climate Research Division, Environment and Climate Change Canada, Toronto, Canada

⁵CNRM, Université de Toulouse, Météo-France, CNRS, Toulouse, France

⁶Met Office Hadley Centre, FitzRoy Road, Exeter, United Kingdom

⁷Université de Lorraine, AgroParisTech, INRAE, UMR Silva, Nancy, France

⁸School of Atmospheric Sciences, Sun Yat-sen University, Guangzhou, China

⁹Instituto Dom Luiz, Faculdade de Ciências, Universidade de Lisboa, Lisbon, Portugal

¹⁰Centre for Hydrology, University of Saskatchewan, Saskatoon, Canada

¹¹WSL Institute for Snow and Avalanche Research SLF, Davos, Switzerland

¹²Institute of Water Problems, Russian Academy of Sciences, Moscow, Russia

¹³Institut für Küstenforschung, Helmholtz-Zentrum Geesthacht, Geesthacht, Germany

1

Early Online Release: This preliminary version has been accepted for publication in *Bulletin of the American Meteorological Society*, may be fully cited, and has been assigned DOI 10.1175/BAMS-D-19-0329.1. The final typeset copyedited article will replace the EOR at the above DOI when it is published.

- ¹⁴CSIRO Oceans and Atmosphere, Canberra, ACT, Australia
- ¹⁵Institute of Industrial Science, the University of Tokyo, Tokyo, Japan
- ¹⁶Grenoble Alpes, Université de Toulouse, Météo-France, CNRS, CNRM, Centre d'Etudes de la Neige, Grenoble, France
- ¹⁷Department of Geography, University of Innsbruck, Innsbruck, Austria
- ¹⁸Meteorological Research Institute, Japan Meteorological Agency, Tsukuba, Japan
- ¹⁹Institute of Meteorology and Climate Research, Karlsruhe Institute of Technology, Karlsruhe, Germany
- ²⁰A.M. Obukhov Institute of Atmospheric Physics, Russian Academy of Sciences, Moscow, Russia
- ²¹Cooperative Institute for Research in Environmental Science/Earth System Research Laboratory, NOAA, Boulder, CO, USA
- ²²Advanced Study Program, National Center for Atmospheric Research, Boulder, Colorado
- ²³Institute of Geography, Russian Academy of Sciences, Moscow, Russia
- ²⁴Department of Atmospheric and Oceanic Sciences, University of Colorado Boulder, Boulder, CO, USA

Abstract

1 Twenty-seven models participated in the Earth System Model - Snow Model Intercomparison
2 Project (ESM-SnowMIP), the most data-rich MIP dedicated to snow modelling. Our findings
3 do not support the hypothesis advanced by previous snow MIPs: evaluating models against
4 more variables, and providing evaluation datasets extended temporally and spatially does not
5 facilitate identification of key new processes requiring improvement to model snow mass and
6 energy budgets, even at point scales. In fact, the same modelling issues identified by previous
7 snow MIPs arose: albedo is a major source of uncertainty, surface exchange parametrizations
8 are problematic and individual model performance is inconsistent. This lack of progress is
9 attributed partly to the large number of human errors that led to anomalous model behaviour
10 and to numerous resubmissions. It is unclear how widespread such errors are in our field and
11 others; dedicated time and resources will be needed to tackle this issue to prevent highly
12 sophisticated models and their research outputs from being vulnerable because of avoidable
13 human mistakes. The design of and the data available to successive snow MIPs were also
14 questioned. Evaluation of models against bulk snow properties was found to be sufficient for
15 some but inappropriate for more complex snow models whose skills at simulating internal
16 snow properties remained untested. Discussions between the authors of this paper on the
17 purpose of MIPs revealed varied, and sometimes contradictory, motivations behind their
18 participation. These findings started a collaborative effort to adapt future snow MIPs to
19 respond to the diverse needs of the community.

Capsule

The latest snow model intercomparison identified the same modelling issues as previous iterations over 23 years. Lack of new insights are attributed partly to human errors and intercomparison projects design.

1 1. Introduction

2

3 The Earth System Model-Snow Model Intercomparison Project (ESM-SnowMIP; Krinner et
4 al., 2018) is the third in a series of MIPs spanning seventeen years investigating the
5 performance of snow models. It is closely aligned with the Land Surface, Snow and Soil
6 Moisture Model Intercomparison Project (LS3MIP; van den Hurk et al. 2016), which is a
7 contribution to the sixth Coupled Model Intercomparison Project (CMIP6). The Tier 1
8 reference site simulations (Ref-Site in Krinner et al., 2018), the results of which are discussed
9 in this paper, is the first of ten planned ESM-SnowMIP experiments and the latest iteration of
10 MIPs using in situ data for snow model evaluation. The Project for Intercomparison of Land
11 surface Parameterization Schemes Phase 2(d) (PILPS 2(d)) was the first comprehensive
12 intercomparison focusing on the representation of snow in land surface schemes (Pitman and
13 Henderson-Sellers, 1998; Slater et al., 2001) and evaluated models at one open site for 18
14 years. It was followed by the first SnowMIP (hereafter SnowMIP1; Etchevers et al., 2002;
15 Etchevers et al., 2004), which evaluated models at four open sites for a total of 19 site-years
16 and by SnowMIP2 (Rutter et al., 2009; Essery et al., 2009) which investigated simulations at
17 five open and forested site pairs for 9 site-years.

18 Twenty-seven models from twenty-two modelling teams participated in the ESM-
19 SnowMIP Ref-Site experiment (ESM-SnowMIP hereafter). A short history of critical findings in
20 previous MIPs is necessary to contextualise the results. PILPS 2(d) identified sources of model
21 scatter to be albedo and fractional snow cover parametrizations controlling the energy
22 available for melt, and longwave radiative feedbacks controlled by exchange coefficients for
23 sensible and latent heat fluxes in stable conditions (Slater et al., 2001). SnowMIP1

24 corroborated the latter finding, adding that the more complex models were better able to
25 simulate net longwave radiation but both complex models and simple models with
26 appropriate parametrizations were able to simulate albedo well (Etchevers et al, 2004) (
27 Baartman et al., 2020, showed that there is no general consensus about what “model
28 complexity” is; for clarity, we will define models explicitly incorporating larger numbers of
29 processes, interactions and feedbacks as more complex). SnowMIP2 found little consistency
30 in model performance between years or sites and, as a result, there was no subset of better
31 models (Rutter et al., 2011). The largest errors in mass and energy balances were attributed
32 to uncertainties in site-specific parameter selection rather than to model structure. All these
33 projects concluded that more temporal and spatial data would improve our understanding of
34 snow models and reduce the uncertainty associated with process representations and
35 feedbacks on the climate.

36 This paper discusses results from model simulations at five mountain sites (Col de Porte,
37 France; Reynolds Mountain East, Idaho, USA; Senator Beck and Swamp Angel, Colorado, USA;
38 Weissfluhjoch, Switzerland), one urban maritime site (Sapporo, Japan) and one Arctic site
39 (Sodankylä, Finland); results for three forested sites will be discussed in a separate
40 publication. Details of the sites, forcing and evaluation data are presented in Menard et al.
41 (2019). Although the 97 site-years of data for these seven reference sites may still be
42 insufficient, they do respond to the demands of previous MIPs by providing more sites in
43 different snowy environments over more years.

44

45 2. The false hypothesis

46 In fiction, a false protagonist is one who is presented as the main character but turns out
47 not to be, often by being killed off early (e.g. Marion Crane in Psycho, 1960; Dallas in Alien,
48 1979; Ned Stark in A Game of Thrones, Martin, 1996). This narrative technique is not used in
49 scientific literature, even though many scientific hypotheses advanced in project proposals
50 are killed early at the research stage. Most scientific journals impose strict manuscript
51 composition guidelines to encourage research studies to be presented in a linear and cohesive
52 manner. As a consequence, many “killed” hypotheses are never presented, and neither are
53 the intermediary steps that lead to the final hypothesis. This is an artifice that we all comply
54 with even though hypothesizing after the results are known (known as HARKing; Kerr, 1998)
55 is a practice associated with the reproduction crisis (Munafò et al., 2017).

56 Our working hypothesis was formed at the design stage of ESM-SnowMIP and is explicit
57 in Krinner et al. (2018): more sites over more years will help us to identify crucial processes
58 and characteristics that need to be improved as well as previously unrecognized weaknesses
59 in snow models. However, months of analysing results led us to conclude the unexpected:
60 more sites, more years and more variables do not provide more insight into key snow
61 processes. Instead, this leads to the same conclusions as previous MIPs: albedo is still a major
62 source of uncertainty, surface exchange parametrizations are still problematic, and individual
63 model performance is inconsistent. In fact, models are less classifiable with results from more
64 sites, years and evaluation variables. Our initial, or false, hypothesis had to be killed off.

65 Developments *have* been made, particularly in terms of the complexity of snow process
66 representations, and conclusions from PILPS2(d) and snow MIPs have undoubtedly driven
67 model development. Table 1 shows that few participating models now have a fixed snow

68 density or thermal conductivity, only two models still parametrize snow albedo as a simple
69 function of temperature, no model uses constant surface exchange coefficients, more models
70 can now represent liquid water in snow, and only three still have a composite snow/soil layer.
71 These changes demonstrate progress for individual models, but they do not for snow science:
72 most of these parametrizations have existed for decades. Differences between models
73 remain, but the range of model complexity is smaller than it was in previous MIPs.

74 The pace of advances in snow modelling and other fields in climate research is limited by
75 the time it takes to collect long-term datasets and to develop methods for measuring complex
76 processes. Furthermore, the logistical challenges of collecting reliable data in environments
77 where unattended instruments are prone to failure continue to restrict the spatial coverage
78 of quality snow datasets.

79 False protagonists allow narrators to change the focus of the story. Our “false hypothesis”
80 allows us to re-focus our paper not on what the model results are – doing so would merely
81 repeat what previous snow MIPs have concluded – but on why, in the twenty four years since
82 the start of PILPS 2 (d), the same modelling issues have repeatedly limited progress in our
83 field, when other fields relying on technology and computing have changed beyond
84 recognition.

85

86 3. The Beauty Contest

87

88 Ranking models (or the “beauty contest”, as insightfully described by Ann Henderson-
89 Sellers when presenting results from PILPS) offers little or no insight into their performance,

90 but it has become the compulsory starting point for presenting MIP results. Figures 1 and 2
91 show models ranked according to errors in daily averages of snow water equivalent (SWE),
92 surface temperature, albedo and soil temperature (note that not all of these variables were
93 measured at all sites or output by all models). To avoid errors in simulations for snow-free or
94 partially snow-covered ground, errors in albedo and surface and soil temperatures were only
95 calculated for periods with measured snow depths greater than 0.1 m and air temperatures
96 below 0°C. Measured and modelled snow surface temperatures greater than 0°C and albedos
97 less than 0.5 were excluded from the error calculations. Bias is shown for SWE, surface
98 temperature, albedo and soil temperature. Root mean square error normalised by standard
99 deviation (NRMSE) is presented only for SWE and surface temperature because standard
100 deviations of albedo and soil temperature are small during periods of continuous snow cover.

101 Discussion of the results in Sections 3.1 to 3.3. will demonstrate why our initial hypothesis
102 was rejected: no patterns emerge, no sweeping statements can be made. The preliminary
103 conclusion presented in Krinner et al. (2018) that “model complexity per se does not explain
104 the spread in performance” still stands. For example, Table 1 shows that RUC is one of the
105 simplest models, but Figures 1 and 2 show that it often has smaller errors than more complex
106 models. This is not to say that model developments are useless: there are large differences
107 between simulations submitted for older and newer versions of a few models. Errors in SWE
108 – the most commonly used variable for evaluation of site simulations – are greatly reduced in
109 HTESSEL-ML, JULES-UKESM/JULES-GL7 and ORCHIDEE-E/ORCHIDEE-MICT compared with
110 HTESSEL, JULES-I and ORCHIDEE-I, and errors in soil temperature are greatly reduced in
111 JSBACH-PF which, unlike its predecessor JSBACH, includes a soil freezing parametrization.
112 There is little or no reduction in errors for other variables between versions.

113 Errors in the ESM-SnowMIP driving and evaluation data are not discussed here because
114 they are discussed in Menard et al. (2019): implicit in the following sections is that a model
115 can only be as good as the data driving it and against which it is evaluated.

116

117 3.1 Snow water equivalent and surface temperature

118

119 Mean SWE and surface temperature NRMSEs in Figure 1 are generally low: below 0.6
120 for half of the models and 1 or greater for only four models. Biases are also relatively low: less
121 than 2°C in surface temperature and less than 0.2 in normalised SWE for four out of five sites
122 in Figure 2. The sign of the biases in surface temperature are the same for at least four out of
123 five sites for all except four models (JULES-I, ORCHIDEE-E, ORCHIDEE-MICT and SWAP). The
124 six models with the largest negative biases in SWE are among the seven models that do not
125 represent liquid water in snow. The seventh model, RUC, has its largest negative bias at
126 Sapporo, where rain-on-snow events are common. Wind-induced snow redistribution, which
127 no model simulates at a point, is partly responsible for Senator Beck being one of the two
128 sites with largest SWE NRMSE in more than half of the models.

129 Four of the best models for SWE NRMSE are among the worst for surface temperature
130 NRMSE (SPONSOR, Crocus, CLASS and HTESSEL-ML). Decoupling of the snow surface from the
131 atmosphere under stable conditions is a long-standing issue which Slater et al. (2001)
132 investigated in PILPS 2(d). Underestimating snow surface temperature leads to a colder
133 snowpack that takes longer to melt and remains on the ground for longer. In 2001, most
134 models used Richardson numbers to calculate surface exchange; in 2019, most use Monin-
135 Obukhov similarity theory (MOST). However, assumptions of flat and horizontally

136 homogeneous surfaces and steady-state conditions in MOST make it inappropriate for
137 describing conditions not only over snow surfaces, but also over forest clearings and
138 mountains: in other words, at all sites in this study. Exchange coefficient are commonly used
139 to tune near-surface temperature in numerical weather prediction models even if to the
140 detriment of the representation of stable boundary layers (Sandu et al., 2013). Conway et al.
141 (2018) showed that such tuning in snowpack modelling improved surface temperature
142 simulations but at the expense of overestimating melt. It is beyond the scope of this paper
143 (and in view of the discussion on sources of errors in Section 4, possibly beyond individual
144 modelling teams) to assess how individual models have developed and evaluated their
145 surface exchange and snowpack evolution schemes. However, differences in model ranking
146 between SWE and surface temperature suggest that this issue is widespread and warrants
147 further attention.

148

149 3.2 Albedo

150

151 Errors in modelled winter albedo (Li et al., 2016) and implications for snow albedo
152 feedback on air temperature (Randall et al., 2007; Flato et al., 2013) have been linked to errors
153 in snow cover fraction (SCF) (e.g. Roesch et al, 2006) and vegetation characteristics in the
154 boreal regions, rather than to the choice or complexity of snow albedo schemes (Essery, 2013;
155 Wang et al, 2016). These should not affect ESM-SnowMIP because vegetation characteristics
156 were provided to participants (all sites discussed here are in clearings or open landscapes)
157 and snow cover during accumulation is expected to be complete. However, eleven models
158 did not impose complete snow cover (Figure 3) such that, again, differences in surface albedo

159 are inextricably linked to differences in snow cover fraction; implications are discussed in
160 Section 4.1.

161 As in previous studies (e.g. Etchevers et al., 2004; Essery, 2013), the specific albedo
162 scheme or its complexity does not determine model performance in ESM-SnowMIP. Neither
163 of the two models with the smallest range of biases, CLASS and EC-Earth, imposed SCF = 1
164 and both use simple albedo schemes in which snow albedo decreases depending on time and
165 temperature. Snow albedo parametrizations (Table 1) determine rates at which albedo varies,
166 but ranges within which the schemes operate are still determined by user-defined minimum
167 and maximum snow albedos to which models are very sensitive. For most models these
168 parameters are the same at all sites, but measurements suggest that they differ; it is unclear
169 whether some of these variations are due to site-specific measurement errors (e.g.
170 instruments or vegetation in the radiometer field of view). This issue should be investigated
171 further as this is not the first time that model results have been inconclusive because of such
172 uncertainties (e.g. Essery et al., 2013).

173

174 3.3 Soil temperature

175

176 Five models systematically underestimate soil temperatures under snow (JSBACH
177 MATSIRO, ORCHIDEE-I, RUC and SURFEX-ISBA) and four systematically overestimate them
178 (CLM5, CoLM, JULES-GL7 and ORCHIDEE-MICT), although negative biases are often larger than
179 positive ones. Soil temperatures are not consistently over- or underestimated by all models
180 at any particular site. Three of the models (JSBACH, JULES-I and ORCHIDEE-I) still include a
181 thermally composite snow-soil layer, and the lack of a soil moisture freezing representation

182 in JSBACH causes soil temperatures to be underestimated. Although newer versions of these
183 models (ORCHIDEE-E, ORCHIDEE-MICT, JSBACH-PF, JULES-GL7 and JULES-UKESM) include
184 more realistic snow-soil process representations, cold biases of the implicit versions have,
185 with the exception of ORCHIDEE-E, been replaced by warm biases, and of similar magnitude
186 between JULES-I and JULES-GL7.

187

188 4. Discussion

189

190 4.1 Motivation behind participation

191

192 One of the motivations behind the design of ESM-SnowMIP was to run a stand-alone
193 MIP dedicated to snow processes parallel to other MIPs, most notably CMIP6 and LS3MIP:
194 “Combining the evaluation of these global-scale simulations with the detailed process-based
195 assessment at the site scale provides an opportunity for substantial progress in the
196 representation of snow, particularly in Earth system models that have not been evaluated in
197 detail with respect to their snow parameterizations” (Krinner et al., 2018). Identifying errors
198 in ESM-SnowMIP site simulations could be linked to model processes that also operate in
199 LS3MIP global simulations, separately from meteorological and ancillary data errors.
200 However, LS3MIP and ESM-SnowMIP results are not directly comparable because land
201 surface schemes (LSSs) include parametrizations that describe sub-grid heterogeneity and
202 some LSSs allow them to be switched off or modified for point simulations. Tables 1 and 2
203 show whether models participated in both MIPs and whether they used point simulation-
204 specific snow cover parametrizations, which is critical for albedo and the most common

205 parametrization to simulate sub-grid heterogeneity. Of the eleven models that did not adjust
206 their sub-grid parametrizations or impose complete snow cover (Figure 3), only one (CLASS)
207 is not participating in LS3MIP. Of those that are participating, three switched off their sub-
208 grid parametrizations (MATSIRO, RUC, and SURFEX-ISBA). Had it been anticipated at the
209 design stage that some models would have considered ESM-SnowMIP to be a means to
210 evaluate their LS3MIP set-up against in situ data, ESM-SnowMIP instructions would have
211 advised to switch off all sub-grid processes; treating a point simulation like a spatial simulation
212 makes evaluating some variables against point measurements futile. This is best illustrated
213 with ORCHIDEE, the three versions of which have the highest negative albedo biases; not only
214 was complete snow cover not imposed, but also the maximum albedo for deep snow on grass
215 (i.e. 0.65 at all sites except Weissfluhjoch) accounts implicitly for sub-grid heterogeneity in
216 large-scale simulations.

217 Although called ESM-SnowMIP, the site simulations were always intended to include
218 physically based snow models that are not part of an ESM but have other applications (Krinner
219 et al., 2018). Table 3 lists what motivated different groups to participate in ESM-SnowMIP
220 Although not explicit in Table 3 because of the anonymity of the comments, for developers of
221 snow physics models, the motivation to participate in a MIP dedicated to scrutinizing the
222 processes they investigate is self-evident. On the other hand, most land surface schemes were
223 first developed to provide the lower boundary conditions to atmospheric models. Because of
224 the dramatic differences in the energy budget of snow-free and snow-covered land, the main
225 requirement for snow models in some LSSs is still just to inform atmospheric models of
226 whether there is snow on the ground or not. The size of the modelling group also matters;
227 more models supported by a single individual or small teams listed exposure as one of their
228 motivations. This discussion revealed that many participants suffered from the “false

229 consensus effect” (Lee et al., 1977), also observed among geoscientists but not explicitly
230 named by Baartman et al . (2020), i.e. they assumed their motivations were universal, or at
231 the very least, widespread. Ultimately, the prestige of MIPs means that, regardless of
232 workload, personal motivation or model performance, they have become compulsory
233 promotional exercises that we cannot afford not to participate in, for better or worse.

234

235 4.2 Errare humanum est

236

237 The increasing physical complexity of models makes them harder for users to
238 understand. Many LSSs are “community” models (e.g. CLM, CoLM, JULES, SURFEX-ISBA),
239 meaning that they are being developed and used by a broad range of scientists whose
240 research interests, other than all being related to some aspect of the land surface, do not
241 necessarily overlap. In many cases, new parametrizations are added faster than old ones are
242 deprecated, causing ever-growing user interfaces or configuration files to become
243 incomprehensible. Benchmarking should help scientists verify that newer versions of a model
244 can reproduce the same results as older versions, but the lag between scientific
245 improvements (hard code) and those at the user interface (soft code) can cause model errors
246 to be introduced by simple avoidable mistakes. The JULES configuration files, for example,
247 contain approximately 800 switches and parameters. Although GL7 and UKESM are the
248 official JULES configurations implemented in the CMIP6 Physical Model and Earth System
249 setups respectively, the ESM-SnowMIP results had to be re-submitted multiple times because
250 large errors were eventually traced to a poorly documented but highly sensitive parameter.

251 It should be noted that JULES and many other models were not intended for point
252 simulations, increasing the possibility of errors in reconfiguring them for ESM-SnowMIP.

253 A different philosophy from some other MIPs has been followed here such that
254 resubmission of simulations was encouraged if initial results did not appear to be
255 representative of the intended model behaviour. Table 4 provides details of the hard- and
256 soft-coded errors identified as a result of discussions that led to sixteen of the twenty-six
257 models re-submitting their results, some more than once. One model was excluded at a late
258 stage because the modelling team did not identify the source of some very large errors that
259 caused the model to be an outlier in all analyses and, therefore, would not have added any
260 scientific value to this paper.

261 Model errors can be statistically quantified; quantifying human errors is somewhat
262 more challenging. A methodology widespread in high-risk disciplines (e.g. medicine, aviation
263 and nuclear power), the Human Reliability Assessment, may be the closest analogue, but it is
264 a preventative measure. Concerns about reproducibility and traceability have motivated a
265 push for analogous methodologies in the Geosciences (Gil et al., 2016), but most remain
266 retrospective steps to retrace at the paper writing stage.

267 Figure 4 quantifies the differences in the performance of the two variables (SWE and
268 soil temperature) and models most affected by human errors before and after resubmission.
269 For some models (JULES-GL7, JSBACH-PF, HTESSEL-ML), SWE NRMSE before resubmission are
270 up to five times higher than after and soil temperature bias double that of corrected
271 simulations (ORCHIDEE-I). Human errors in models and, as discussed in Menard et al. (2019)
272 for the first ten reference sites in ESM-SnowMIP, in data are inevitable, and this snow MIP
273 shows that they are widespread. The language we use to describe numerical models has

274 distanced them from the fact that they are not, in fact, pure descriptions of physics but rather
275 equations and configuration files written by humans. *Errare humanum est, perseverare*
276 *diabolicum*. Menard et al. (2015) showed that papers already published had used versions of
277 JULES that included bugs affecting turbulent fluxes and causing early snowmelt. There is no
278 requirement for authors to update papers after publication if retrospective enquiries identify
279 some of the published results as erroneous. In view of the many errors identified here, further
280 investigations are required to start understanding how widespread errors in publications are.
281 Whether present in initialisation files or in the source code, these errors impair or slow
282 progress in our understanding of snow modelling because they misrepresent the ability of
283 models to simulate snow mass and energy balances.

284

285 4.3 Model documentation

286

287 As with many other areas of science, calls for reproducibility of model results to
288 become a requirement for publication are gaining ground (Gil et al., 2016). Table 1 was initially
289 intended to list the parametrizations considered most important in snow modelling (Essery
290 et al., 2013; Essery, 2015), with, as is conventional (e.g. Rutter et al., 2009; Krinner et al.,
291 2018), a single reference per model. Referencing the parametrizations in the twenty-seven
292 models requires, in fact, seventy-nine papers and technical reports; a more detailed version
293 of the table and associated references are included in the supplementary material. The lead-
294 author first identified fifty-one references, and the modelling teams then provided references
295 to fill the remaining gaps. However, some suggested the wrong references, others revised
296 their initial answers and a few even discovered that some parametrizations are not described

297 at all. Not only is it extremely rare to find complete documentation of a model in a single
298 publication, it is also difficult to find all parametrizations described at all in the literature.
299 When this happens, some parametrizations are described in publications for other models.
300 Often, the most recent publication refers to previous ones, which may or may not be the first
301 to have described the model, comprehensively or not. Incomplete documentation would be
302 an annoying but unimportant issue if this exercise had not led to the identification of some of
303 the errors discussed in Section 4.2.

304 Less than a decade ago, it was at best difficult and at worst impossible to publish scientific
305 model descriptions. The open access culture, issues of reproducibility and online platforms
306 dedicated to publication of source code and data have reversed this trend such that it is now
307 difficult to imagine research relying on a new model with proprietary code being published.
308 Yet, it is a truth universally acknowledged that openly budgeting in a project proposal for the
309 added time it takes to publish comprehensive data and model descriptions is unadvisable,
310 despite many funding bodies enforcing open-access policies. The problem remains for models
311 developed before the tide changed. Two examples illustrate this best. The first concerns the
312 number of papers which refer to Anderson (1976) for snow density, liquid water retention or
313 thermal conductivity. Equations for these parametrizations do appear in the report, but often
314 not in the form presented in subsequent papers (Essery et al., 2012 pointed out that most
315 actually use the forms in Jordan, 1991), or they are themselves reproductions of equations
316 from earlier studies (especially for snow thermal conductivity). The second example is a quote
317 taken from the paper describing VEG3D (Braun and Schädler, 2005): “The snow model is
318 based on the Canadian Land Surface Scheme (CLASS) (Verseghy 1991) and ISBA (Douville et
319 al. 1995) models, and accounts for changes of albedo and emissivity as well as processes like
320 compaction, destructive metamorphosis, the melting of snow, and the freezing of liquid

321 water.” This sentence is the only description in English of the snow model in VEG3D; a more
322 comprehensive description, not referenced in Braun and Shädler (2005), is available in
323 German in a PhD thesis (Grabe, 2002). The study in which the quote appears did not focus on
324 snow processes, so a full description of the snow model may not have been necessary, but it
325 is nonetheless a cause for concern that referees, at the very least, did not require clarifications
326 as to which processes were based on CLASS and which on ISBA. Changes in emissivity certainly
327 were not based on either model as both did – and still do – have fixed emissivity. This is the
328 most succinct description of a snow model, but not the only one to offer little or no
329 information about process representations. At the other end of the spectrum, the CLM5
330 documentation is the most comprehensive and makes all the information available in a single
331 technical report (Lawrence et al., 2020). A few models follow closely with most information
332 being available in a single document that clearly references where to obtain additional
333 information (e.g. CLASS, SURFEX-ISBA, HTESSEL, JULES, SNOWPACK). The “Publish or perish”
334 culture is estimated to foster a nine percent yearly growth rate in scientific publications
335 (Bornmann and Mutz, 2015) which will be matched by a comparable rate of solicitations for
336 peer reviewing. Whether it is because we do not take or have time to fact-check references,
337 the current peer-review process is failing when poorly described models are published. The
338 aim of LS3MIP and ESM-SnowMIP is to investigate systematic errors in models; errors can be
339 quantified against evaluation data for any model, but poor documentation accentuates our
340 poor understanding of model behaviour and reduces MIPs to statistical exercises rather than
341 to insightful studies.

342

343 5. What the future holds

344

345 Historically, PILPS (Henderson-Sellers et al., 1995) and other intercomparison projects
346 have provided platforms to motivate model developments; they are now inextricably linked
347 to successive IPCC reports. In view of heavily mediatised errors such as the claim that
348 Himalayan glaciers would melt by 2035 – interestingly described as “human error” by the then
349 IPCC chairman Rajendra Pachauri (archive.ipcc.ch, 2010; Times of India, 2010) – we must
350 reflect on how damaging potential errors are to the climate science community. Not only are
351 the IPCC reports the most authoritative in international climate change policy-making, but
352 they have become – for better or worse – proxies for the credibility of climate scientists to
353 the general public. It is therefore time that we reflect on our community and openly
354 acknowledge that some model uncertainties cannot be quantified at present because they
355 are due to human errors.

356 Other factors are also responsible for the modelling of snow processes not having
357 progressed as fast as other areas relying on technology. Discussions on the future of snow
358 MIPs involving organisers and participants of ESM-SnowMIP issued from this study. As in the
359 discussion about motivation of participants, suggestions for the design of future MIPs were
360 varied, and at times contradictory, but responses from participants reflected the purpose
361 their models serve (Table 4). The IPCC Expert Meeting on Multi Model Evaluation Good
362 Practice Guidance states that “there should be no minimum performance criteria for entry
363 into the CMIP multi-model database. Researchers may select a subset of models for a
364 particular analysis but should document the reasons why” (Knutti et al., 2010). Nevertheless,
365 many participants argued that the “one size fits all” approach should be reconsidered. ESM-

366 SnowMIP evaluated models against the same bulk snowpack properties as previous snow
367 MIPs. This suited LSSs that represent snow as a composite snow/soil layer or as a single layer,
368 but there is a demand for more complex models that simulate profiles of internal snowpack
369 properties to be evaluated against data that match the scale of the processes they represent
370 (e.g. snow layer temperatures, liquid water content and microstructure). Models used at very
371 high resolution for avalanche risk forecasting (such as Crocus and SNOWPACK; Morin et al.,
372 2020) and by the tourism industry are constantly being tested during the snow season and
373 errors can cost lives and money. However, obtaining reliable data and designing appropriate
374 evaluation methodologies to drive progress in complex snow models is challenging (Menard
375 et al., 2019). For example, solving the trade-off between SWE and surface temperature errors
376 requires more measurements of surface mass and energy balance components: simple in
377 theory but expensive and logistically difficult in practice. The scale at which even the more
378 complex models operate is also impeding progress. Until every process can be described
379 explicitly, the reliance of models on parametrizations to describe very small scale processes
380 (such as the surface exchanges upon which the above trade-off depends) are inevitable
381 sources of uncertainty.

382 Despite expressing a need for change in the design of snow MIPs, many participants
383 described ESM-SnowMIP as a success because it allowed them to identify bugs or areas of
384 their models in need of further improvements; some improvements were implemented in the
385 course of this study, others are in development. Ultimately, ESM-SnowMIP's main flaw is of
386 not being greater than the sum of its parts. Its working hypothesis was not supported and,
387 per se, has failed to advance our understanding of snow processes. However, the
388 collaborative effort allowed us to report a false, but plausible hypothesis, to expose our
389 misplaced assumptions and to reveal a disparity of opinions on the purpose, design and future

390 of snow MIPs. In view of our findings, of the time investment required of participating
391 modellers and of novel ways to utilise already available global-scale simulations (e.g. Mudryk
392 et al., 2020), most planned ESM-SnowMIP experiments may not go ahead, but site simulations
393 with evaluation data covering bulk and internal snowpack properties will be expanded.
394 Learning from our mistakes to implement future MIPs may yet make it an unqualified success
395 in the long term.

396

397 Acknowledgments

398

399 CM and RE were supported by NERC grant NE/P011926/1. Simulations by participating
400 models were supported by the following programs and grants: Capability Development Fund
401 of CSIRO Oceans and Atmosphere, Australia (CABLE); Canada Research Chairs and Global
402 Water Futures (CRHM); H2020 APPLICATE grant 727862 (HTESSEL); Met Office Hadley Centre
403 Climate Programme by BEIS and Defra (JULES-UKESM and GL7); TOUGOU Program from
404 MEXT, Japan (MATSIRO); RUC by NOAA grant NOAA/NA17OAR4320101 (RUC); Russian
405 Foundation for Basic Research grant 18-05-60216 (SPONSOR); Russian Science Foundation
406 Grant 16-17-10039 (SWAP). ESM-SnowMIP was supported by the World Climate Research
407 Programme's Climate and Cryosphere (CliC) core project.

408

409

410

411

412 References

413

414 Abramowitz, G. and Bishop, C. H., 2015: Climate model dependence and the ensemble
415 dependence transformation of CMIP projections, *Journal of Climate*, 28, 2332–2348.

416 Alien, 1979. [film]. Directed by R. Scott.

417 Anderson, E., 1976: A point energy and mass balance model of a snow cover. NOAA Tech
418 Rep. NWS 19, 150 pp.

419 Archive.ipcc.ch, IPCC statement on the melting of Himalayan glaciers
420 <https://archive.ipcc.ch/pdf/presentations/himalaya-statement-20january2010.pdf>, Accessed
421 on 29 October 2019.

422 Baartman, J.E.M, Melsen, L.A., Moore, D. and van der Ploeg, M.J., 2020: On the complexity
423 of model complexity: Viewpoints across the geosciences, *Catena*, 186,
424 10.1016/j.catena.2019.104261, <https://doi.org/10.1016/j.catena.2019.104261>

425 Bornmann, L. and Mutz, R., 2015: Growth rates of modern science: A bibliometric analysis
426 based on the number of publications and cited references, *Journal of the Association for*
427 *Information Science and Technology*, 66, 2215-2222, <https://doi.org/10.1002/asi.23329>

428 Braun, F. J. and Schädler, G., 2005: Comparison of Soil Hydraulic Parameterizations for
429 Mesoscale Meteorological Models, *J. Appl. Meteorol.*, 44, 1116–1132,
430 <https://doi.org/10.1175/JAM2259.1>.

431 Clark, M. P. and Coauthors, 2015: A unified approach for process-based hydrologic modeling:
432 1. Modeling concept, *Water Resour. Res.*, 51, 2498– 2514, doi:10.1002/2015WR017198.

433 Conway J.P., Pomeroy J.W., Helgason W.D. and Kinar, N.J., 2018: Challenges in modelling
434 turbulent heat fluxes to snowpacks in forest clearings, *Journal of Hydrometeorology*, 19,
435 1599–1616, <https://doi.org/10.1175/JHM-D-18-0050.1>.

436 Douville, H. and Mahfouf, J.F., 1995: A new snow parameterization for the Météo-France
437 Climate Model. 1. Validation in stand-alone experiments. *Climate Dynamics*, 12, 21–35.

438 Essery, R., 2013: Large-scale simulations of snow albedo masking by forests, *Geophys. Res.*
439 *Lett.*, 40, 5521– 5525, doi:10.1002/grl.51008.

440 Essery, R., 2015: A Factorial snowpack model (FSM 1.0) *Geoscientific Model Developemnt*, 8,
441 3867-3876, <https://doi.org/10.5194/gmd-8-3867-2015>

442 Essery, R. and Etchevers, P., 2004: Parameter sensitivity in simulations of snowmelt, *Journal*
443 *of Geophysical Research*, 109, 1-15, <https://doi.org/10.1029/2004JD005036>.

444 Essery, R., Morin, S., Lejeune, Y. and Menard, C., 2013: A comparison of 1701 snow models
445 using observations from an alpine site. In: *Advances in Water Resources*, 55, 131-148.

446 Essery, R., Rutter, N., Pomeroy, J., Baxter, R., Stähli, M., Gustafsson, D., Barr, A., Bartlett, P.,
447 and Elder, K., 2009: SNOWMIP2: An Evaluation of Forest Snow Process Simulations, *B. Am.*
448 *Meteorol. Soc.*, 90, 1120–1135, <https://doi.org/10.1175/2009BAMS2629.1>.

449 Etchevers, P. and Coauthors, 2002: SnowMiP, an intercomparison of snow models: first
450 results. In: *Proceedings of the International snow science workshop*, Penticton, Canada, 29
451 Sep.-4 Oct. 2002.

452 Etchevers, P. and Coauthors, 2004: Validation of the surface energy budget simulated by
453 several snow models (SnowMIP project), *Annals of Glaciology*, 38, 150-158.

454 Flato, G and Coauthors, 2013: Evaluation of Climate Models. In: Climate Change 2013: The
455 Physical Science Basis. Contribution of Working Group I to the Fifth Assessment Report of the
456 Intergovernmental Panel on Climate Change [Stocker, T.F., D. Qin, G.-K. Plattner, M. Tignor,
457 S.K. Allen, J. Boschung, A. Nauels, Y. Xia, V. Bex and P.M. Midgley (eds.)]. Cambridge University
458 Press, Cambridge, United Kingdom and New York, NY, USA.

459 Gil, Y. and Coauthors, 2016: Toward the Geoscience Paper of the Future: Best practices for
460 documenting and sharing research from data to software to provenance, *Earth and Space*
461 *Science*, 3, 388-415, <https://doi.org/10.1002/2015EA000136>.

462 Grabe, F., 2002: Simulation der Wechselwirkung zwischen Atmosphäre, Vegetation und
463 Erdoberfläche bei Verwendung unterschiedlicher Parametrisierungsansätze (PhD thesis,
464 Karlsruhe University, Karlsruhe). Retrieved from KIT-Bibliothek (KITopen-ID: 122002).

465 Henderson-Sellers, A., Pitman, A.J., Love, P.K., Irannejad, P. and Chen, T., 1995: The project
466 for Intercomparison of land surface parameterisation schemes (PILPS) Phases 2 and 3. *Bull.*
467 *Amer. Meteor. Soc.*, 76, 489-503.

468 Jordan, R., 1991: A One-Dimensional Temperature Model for a Snow Cover, Technical
469 Documentation for SNTHERM.89, Special Report 91-16, U.S. Army Corps of Engineers, 62 pp.

470 Kerr, N. L., 1998: HARKing: hypothesizing after the results are known. *Personality and Social*
471 *Psychology Review*, 2, 196–217.

472 Knutti, R., Abramowitz, G., Collins, M., Eyring, V., Gleckler, P.J., Hewitson, B. and Mearns L.,
473 2010: Good Practice Guidance Paper on Assessing and Combining Multi Model Climate
474 Projections. In: Meeting Report of the Intergovernmental Panel on Climate Change Expert
475 Meeting on Assessing and Combining Multi Model Climate Projections [Stocker, T.F., D. Qin,

476 G.-K. Plattner, M. Tignor, and P.M. Midgley (eds.)). IPCC Working Group I Technical Support
477 Unit, University of Bern, Bern, Switzerland.

478 Knutti, R., Sedláček, J, Sanderson, B.M., Lorenz, R., Fischer, E.M. and Eyring, V, 2017: A climate
479 model projection weighting scheme accounting for performance and interdependence.
480 *Geophysical Research Letters*, 44, 1909–1918.

481 Krinner, G. and Coauthors, 2018: ESM-SnowMIP: Assessing models and quantifying snow-
482 related climate feedbacks, *Geosci. Model Dev. Discuss.*, [https://doi.org/10.5194/gmd-2018-](https://doi.org/10.5194/gmd-2018-153)
483 153.

484 Lafaysse, M., Cluzet, B., Dumont, M., Lejeune, Y., Vionnet, V., and Morin, S., 2017: A
485 multiphysical ensemble system of numerical snow modelling, *The Cryosphere*, 11, 1173-1198,
486 <https://doi.org/10.5194/tc-11-1173-2017>.

487 Lawrence, D. and Coauthors, 2020: Technical Description of version 5.0 of the Community
488 Land Model (CLM), NCAR, 329pp,
489 http://www.cesm.ucar.edu/models/cesm2/land/CLM50_Tech_Note.pdf

490 Lee, R., Greene, D. and House, P., 1977: The 'false consensus effect': An egocentric bias in
491 social perception and attribution processes, *Journal of Experimental Social Psychology*, 13,
492 279–301. [https://doi.org/10.1016/0022-1031\(77\)90049-X](https://doi.org/10.1016/0022-1031(77)90049-X).

493 Li, Y., Wang, T., Zeng, Z., Peng, S., Lian, X and Piao S., 2016: Evaluating biases in simulated
494 land surface albedo from CMIP5 global climate models, *Journal of Geophysical Research:*
495 *Atmospheres*, 121, 6178-6190, <https://doi.org/10.1002/2016JD024774>

496 Martin, G.R.R., 1996: *A Game of Thrones*, Bantam Spectra, USA.

497 Menard, C.B., Ikonen, J., Rautiainen, K., Aurela, M., Arslan, A.N., and Pulliainen, J., 2015:
498 Effects of Meteorological and Ancillary Data, Temporal Averaging, and Evaluation Methods
499 on Model Performance and Uncertainty in a Land Surface Model. *J. Hydrometeor.*, 16, 2559–
500 2576, <https://doi.org/10.1175/JHM-D-15-0013.1>.

501 Menard, C. B. and Coauthors, 2019: Meteorological and evaluation datasets for snow
502 modelling at 10 reference sites: description of in situ and bias-corrected reanalysis data, *Earth*
503 *Syst. Sci. Data*, 11, 865–880, <https://doi.org/10.5194/essd-11-865-2019>.

504 Morin S. and Coauthors, 2020: Application of physical snowpack models in support of
505 operational avalanche hazard forecasting: A status report on current implementations and
506 prospects for the future, *Cold Regions Science and Technology*, 170, 102910,
507 <https://doi.org/10.1016/j.coldregions.2019.102910>

508 Mudryk, L., Santolaria-Otín, M., Krinner, G., Ménégos, M., Derksen, C., Brutel-Vuilmet, C.,
509 Brady, M., and Essery, R., 2020: Historical Northern Hemisphere snow cover trends and
510 projected changes in the CMIP-6 multi-model ensemble, *The Cryosphere Discuss.*,
511 <https://doi.org/10.5194/tc-2019-320>.

512 Munafò, M. and Coauthors, 2017: A manifesto for reproducible science, *Nature Human*
513 *Behaviour*, 1, 0021, <https://doi.org/10.1038/s41562-016-0021>.

514 Niu, G.-Y. and Coauthors, 2011: The community Noah land surface model with
515 multiparameterization options (Noah-MP): 1. Model description and evaluation with local-
516 scale measurements. *J. Geophys. Res.*, 116, D12109, doi: 10.1029/2010JD015139.

517 Pitman, A.J. and Henderson-Sellers, A., 1998: Recent progress and results from the project for
518 the intercomparison of land surface parameterization schemes. *Journal of Hydrology*, 212-
519 213, 128-135, [https://doi.org/10.1016/S0022-1694\(98\)00206-6](https://doi.org/10.1016/S0022-1694(98)00206-6)

520 Psycho, 1960: [film] Directed by A. Hitchcock.

521 Randall, D.A. and Coauthors, 2007: Climate Models and Their Evaluation. In: Climate Change
522 2007: The Physical Science Basis. Contribution of Working Group I to the Fourth Assessment
523 Report of the Intergovernmental Panel on Climate Change [Solomon, S., Qin, D., Manning, M.,
524 Chen, Z., Marquis, M., Averyt, K.B., Tignor, M. and Miller, H.L. (eds.)]. Cambridge University
525 Press, Cambridge, United Kingdom and New York, NY, USA.

526 Roesch, A., 2006: Evaluation of surface albedo and snow cover in AR4 coupled climate models,
527 *J. Geophys. Res.*, 111, D15111, doi:10.1029/2005JD006473.

528 Rutter, N. and Coauthors, 2009: Evaluation of forest snow processes models (SnowMIP2), *J.*
529 *Geophys. Res.*, 114, D06111, <https://doi.org/10.1029/2008JD011063>.

530 Sandu, I., Beljaars, A., Bechtold, P., Mauritsen, T., and Balsamo, G., 2013: Why is it so difficult
531 to represent stably stratified conditions in numerical weather prediction (NWP) models?, *J.*
532 *Adv. Model. Earth Syst.*, 5, 117– 133, doi:10.1002/jame.20013.

533 Slater, A.G. and Coauthors, 2001: The representation of snow in land surface schemes: results
534 from PILPS 2(d). *Journal of Hydrometeorology*, 2, 7–25.

535 [https://timesofindia.indiatimes.com/videos/news/Pachauri-admits-mistake-in-IPCC-](https://timesofindia.indiatimes.com/videos/news/Pachauri-admits-mistake-in-IPCC-report/videoshow/5492814.cms)
536 [report/videoshow/5492814.cms](https://timesofindia.indiatimes.com/videos/news/Pachauri-admits-mistake-in-IPCC-report/videoshow/5492814.cms), 2010, Pachauri admits mistake in IPCC report, accessed on
537 4 October 2019.

538 van den Hurk, B. and Coauthors, 2016: LS3MIP (v1.0) contribution to CMIP6: the Land Surface,
539 Snow and Soil moisture Model Intercomparison Project – aims, setup and expected outcome,
540 *Geosci. Model Dev.*, 9, 2809-2832, <https://doi.org/10.5194/gmd-9-2809-2016>.

541 Versegny, D. L., 1991: Class—A Canadian land surface scheme for GCMS. I. Soil model. *Int. J.*
542 *Climatol.*, 11, 111-133, <https://doi.org/10.1002/joc.3370110202>.

543 Wang, L., Cole, J.N.S., Bartlett, P., Versegny, D., Derksen, C., Brown, R., and von Salzen, K.,
544 2016: Investigating the spread in surface albedo for snow-covered forests in CMIP5 models,
545 *Journal of Geophysical Research: Atmospheres*, 121, 1104– 1119,
546 [doi:10.1002/2015JD023824](https://doi.org/10.1002/2015JD023824).

547

548

549

550

*Table 1: Key characteristics of snow model parametrizations and variables on which they depend, and number of papers per model over which descriptions of the seven parametrizations are spread. Abbreviations and symbols: LWC = Liquid water content, SCF = snow cover fraction (“point” means models used point-specific parametrizations, “grid” means they did not), MC = Mechanical compaction, OL = Obukhov length, PC = Personal communication, Ri_b = bulk Richardson number, * = references provided by personal communication and cannot be traced in the existing literature about this specific model. A more detailed version of this table including full references for parametrizations is available in the supplementary material.*

	Albedo	Conductivity	Density	Turbulent fluxes	LWC	SCF	Snow layering	n Papers
CABLE-SLI	Spectral	Power function	MC	OL	Yes	Point	Single	3
CLASS	Spectral	Quadratic equation	Time	Ri _B	Yes	Grid	Single	2
CLM5	Spectral	Density	MC	OL	Yes	Grid	Multi	1
CoLM	Spectral	Quadratic equation	MC	OL	Yes	Grid	Multi	7*
CRHM	Spectral	Density and humidity	MC	OL	Yes	Point	Multi	4* + PC
Crocus	Spectral	Power function	MC	Ri _B	Yes	Point	Multi	3
EC-EARTH	Time and temperature	Power function	MC	OL	Yes	Grid	Single	3*
ESCIMO	Temperature	None	Time	Empirical	Yes	Point	Single	3*
HTESSEL	Time and temperature	Power function	MC	OL	Yes	Grid	Single	3
HTESSEL (ML)							Multi	3
SURFEX-ISBA	Spectral	Power function	MC	Ri _B	Yes	Point	Multi	2
JSBACH	Spectral	Fixed	Fixed	OL	No	Point	Composite	3*
JSBACH3-PF		Power function	Time				Multi	4*
JULES-GL7 JULES-UKESM	Spectral	Power function	MC	OL	Yes	Point	Multi	2
JULES-I	Temperature	Fixed	Fixed	OL	No.	Point	Composite	1
MATSIRO	Spectral	Fixed	Fixed	OL	No.	Point	Multi	3
ORCHIDEE-E ORCHIDEE-MICT	Time	Quadratic equation	MC	OL	Yes	Grid	Multi	1 + PC
ORCHIDEE-I		Fixed	Fixed		No			
RUC	Time	Fixed	MC	OL	No	Grid	Multi	3 + PC
SMAP	Spectral	Quadratic equation	MC	OL	Yes	Point	Multi	3
SNOWPACK	Statistical	Conductivity model	Empirical	OL	Yes	Point	Multi	5
SPONSOR	Time	Density	MC	OL	Yes	Grid	Multi	2 + PC

SWAP	Density	Density	SWE and snow	OL	Yes	Point	Single	3
VEG3D	Time	Density	Time	OL	No	Point	Single	4*

Table 2: Participating models and modelling teams. ESM-SnowMIP provided vegetation height, soil type and snow-free albedo to the participants; where relevant, these may differ from LS3MIP configurations.

Model	ESM-SnowMIP contact	Model type	Model version	Model configuration	Differences between LS3MIP and ESM-SnowMIP configurations
CABLE-SLI	Matthias Cuntz, Vanessa Haverd	LSS in Access	CABLE revision 4252	CABLE including SLI as described in Haverd and Cuntz (2016). Snow and ice extensions as in Cuntz and Haverd (2018). 12 soil layers.	Did not participate in LS3MIP
CLASS	Paul Bartlett	LSS in CanESM	CLASS 3.6.2	CLASS-CTEM off-line code with CTEM turned off, and using the 2-band snow albedo and associated snow-ageing scheme. Initialization files are available on demand. Other than adjustments to match the site properties (e.g. soil type, vegetation, snow-free albedo) all parameters are the model default values.	Did not participate in LS3MIP
CLM5	Sean Swenson	LSS in CESM	CLM5.0	Standard	No difference.
CoLM	Yongjiu Dai, Hua Yuan	LSS in BNU-ESM and CAS-ESM	CoLM Version 2014	Default	CoLM Version 2005 Many differences including pedotransfer functions of soil hydraulic and thermal parameters, numerical solution of Richards

					equation of soil water content.
CRHM	Xing Fang, John Pomeroy	Hydrological model	CRHM 01/17/18	Adapted from CRHM plot-scale simulation project for coniferous forest and forest clearing sites in Canadian Rocky Mountains detailed in Pomeroy et al. (2012) with modified configuration for soil module allowing simulations for permafrost and seasonal frost.	Did not participate in LS3MIP
Crocus	Matthieu Lafaysse	Snow physics model	Git tag ESM- SnowMIP-Crocus- ESCROC (= commit b57f02d6 4/12/2017)	Crocus : default configuration as defined in Lafaysse et al. (2017), Figure 2. Drift module allowing change of physical properties of near surface snow activated for SNB and WFJ.	Did not participate in LS3MIP
EC-EARTH	Emanuel Dutra	LSS in EC-EARTH	EC-EARTH v3.2.2 revision r4381	Offline "OSM" configuration with prescribed surface albedo and vegetation.	LS3MIP simulation will be done with the latest "frozen" model version for CMIP6, including interactive vegetation and variable surface albedo.
ESCIMO	Thomas Marke, Ulrich Strasser	Snow surface energy balance model	ESCIMO v5 based on ESCIMO v4 with additional functionality described in Marke et al. (2016).	Albedo parameterization as in Cox et al. (1999) Sensible heat equation as in Weber (2008) Empirical density function as in Essery et al. (2013)	Did not participate in LS3MIP
HTESSEL HTESSEL- ML	Gabriele Arduini	LSS of ECMWF operational forecasting system	HTESSEL cycle 43r3	Operational HTESSEL configuration uses the single layer snow scheme from Dutra et al. (2010). The experimental HTESSEL configuration (HTESSELML)	Did not participate in LS3MIP

				uses a multi-layer snow scheme documented in Arduini et al. 2019 (under review in JAMES). Note that the configuration of the multi-layer snow scheme and model cycle used for ESM-SnowMIP runs differ from Arduini et al. (2019).	
SURFEX-ISBA	Bertrand Decharme, Aaron Boone	LSS in CNRM-CM	SURFEX version 8.0 (ISBA and all related schemes including snow are embedded in the SURFEX numerical platform)	As in Decharme et al. (2016) denoted as the "NEW" experiment.	Snow grid-cell fraction doesn't account for vegetation in the 1-dimensional ESM-SnowMIP runs.
JSBACH3 JSBACH3 -PF	Stefan Hagemann	LSS in MPI-ESM	JSBACH3 (Revision 9168, state of 31.07.2017) and JSBACH3-PF (same revision but with improved snow parametrizations inherited from JSBACH4)	Time step: 450s, With YASSO soil model, no dynamic vegetation, no nitrogen, no disturbances and no land use transitions. Orography and LAI do not affect surface roughness. Soil states were initialized from previous global offline simulation using GWSP3 forcing. JSBACH3-PF uses the "permafrost" configuration with enabled soil freezing and thawing, and with related processes based on Ekici et al. (2014).	JSBACH-PF did not participate in LS3MIP JSBACH3: No difference
JULES-I	Cecile Menard, Richard Essery	LSS in HadCM3	JULES 4.8 (Revision 7629)	Zero-layer snow model as described in Best et al. (2011).	Did not participate in LS3MIP
JULES-GL7 JULES-UKESM	Eleanor Burke	LSS in HadGEM3-GC3 and UKESM	JULES 5.3	GL7 and UKESM configurations with site-specific characteristics.	Different fractional snow cover parametrization for plot-scale and distributed simulations.

MATSIRO	Tomoko Nitta, Hyungjun Kim	LSS in MIROC	MATSIRO 6	MATSIRO for offline land simulations. The configuration is the same as the GSWP3 simulations except for subgrid-scale parameterizations (tile scheme, SSNOWD snow cover parameterization and arctic wetland scheme), which are turned off for plot-scale simulations.	All subgrid-scale parameterizations are tuned off for plot-scale simulations.
ORCHIDE E-E ORCHIDE E-I ORCHIDE E-MICT	Claire Brutel-Vuilmet, Gerhard Krinner	LSS in IPSL-CM	ORCHIDEE E and I TRUNK revision 4695; ORCHIDEE MICT 8.7.1 revision 5308	TRUNK is the version of ORCHIDEE that is used in the first CMIP6 runs. We have the implicit snow version (TRUNK-I) which is the older snow that was used in CMIP5 and the explicit snow version (TRUNK-E) that is used in CMIP6 (based on Wang et al., 2013). MICT is the high-latitude version of ORCHIDEE (Guimberteau et al., 2018).	No difference.
RUC	Tatiana Smirnova	LSS in NOAA/NCEP operational forecasting systems	RUC model – WRF 4.0 official release	Standard RUC configuration for offline simulations: 9 levels in soil, 2-layer snow model with separate treatment of snow-covered and snow-free areas for patchy snow.	Subgrid-scale parameterizations for fractional snow cover and surface parameters are turned off for ESM-SnowMIP.
SMAP	Masashi Niwano	Snow physics model	SMAP v4.23rc1	SMAP v4.23rc1	Did not participate in LS3MIP
SNOWPACK	Nander Wever, Charles Fierz	Snow physics model	MeteoIO preprocessing library: revision 2011 from https://models.slf.ch/svn/meteoio/trunk SNOWPACK model:	The standard version of SNOWPACK was used, in default configuration.	Did not participate in LS3MIP

			revision 1480 from https://models.slf.ch/svn/snowpack/branches/dev		
SPONSOR	Dmitry Turkov, Vladimir Semenov	Hydrological model	SPONSOR, ver.2.0	The model was adapted for calculations of spatially distributed landscape characteristics with observed meteorological forcing. The latest version of the snow model is described in Turkov and Sokratov (2016).	No difference
SWAP	Olga Nasonova, Yeugeny Gusev	LSS	As described in Gusev and Nasonova (2003)	As described in Gusev and Nasonova (2003)	Did not participate in LS3MIP
VEG3D	Gerd Schädler	Soil and vegetation model	As described in Braun and Schädler (2005)	Standard configuration: 8 soil layers, time step 300 s.	Did not participate in LS3MIP

Motivation behind participation	Future of snow MIPS
<ul style="list-style-type: none"> • To identify key missing processes. • To cut out the noise from ensemble simulations in order to extract the signal. • To compare how models implement snow processes and, if possible, what are the implications. • To have a detailed analysis of one's own model; doing the model simulations is easier than analysing the results. • To provide new insights into modelling. • To document the current state of the models. • To help modellers understand their and other models better. • To determine the skill of an operational model in offline simulations before starting coupled simulations for weather predictions. • To motivate model improvements. • To participate in the beauty contest (the statistical performance of my not-so-sophisticated model is similar to complex process-based models). • To identify a range of "good enough" - models reflecting the range of process uncertainty. 	<ul style="list-style-type: none"> • Allow re-submission of simulations if errors are identified. • Provide model code and initialisation files as well as model results for transparency. • Move towards a more process-based diagnostic in order to improve parametrizations and not just to tune parameters. • Need new evaluation metrics. • Evaluate against internal snowpack properties (e.g. snow layer thermal conductivity, temperature, density). • Move towards fewer models with multiple hypotheses (e.g. FSM, Essery, 2015; SUMMA, Clark et al., 2015; or Noah-MP, Niu et al., 2011) • Cluster models depending on their complexity. • Not all models should be accepted. There could be minimum requirements in terms of parametrizations (e.g. stability dependent exchange coefficients); outliers from the previous experiment would not be allowed to participate in the next stages; new models should present a proof of energy and moisture conservation in their models.

<ul style="list-style-type: none"> • To make one's model visible to the snow modelling community. • To be part of the snow modelling community. • To evaluate one's model at reference sites across different elevation gradients and climatic settings. • To avoid equifinality problems by evaluating models performance with multiple variables that contribute to and are relevant to snow processes. • To provide benchmarks against which to evaluate models. 	<ul style="list-style-type: none"> • All models should be accepted, but different levels of involvement should be allowed so modelling groups can choose the experiments they want to participate in. • Constrain model sensitivity with observations (e.g. SWE, snow albedo) or fixed variables. • Provide evaluation data at the same time as the forcing data. • Provide fewer sites as initialisation of many sites can be a source of human errors. • Provide more challenging sites (e.g. tundra, wind-blown).
--	---

Table 3: Summary of discussions with ESM-SnowMIP participants about (1) what motivated them to participate and (2) their suggestions about the design of the next snow MIP.

Table 4: Hard and soft coded errors identified by the results analysis team (AT) or modelling team (MT) in the course of this study.

	Unusual model behaviour	Model
Soft-coded errors		
Did not change start time between SNB and SWA (start at 00:00) and other sites (start at 01:00)	Mismatched timestamps (AT)	All models
Initial conditions taken from wrong date	Mismatched timestamps (AT)	CLM5
Specified site-specific parameters not taken from site descriptions	Unrealistically low albedo with consequences on snow mass and melt (AT)	JSBACH, JSBACH-PF, JULES-I
Wrong forcing file used for one site	Models results were identical at two sites (AT)	RUC
Simulations used UTC times instead of local times	Unrealistically high albedo (AT)	Crocus
Many variations in output file formats; wrong variable name; variations in the interpretation of the ESM-SnowMIP definition of output variables; different sign conventions.	N/A	One or more in adjacent list for most, if not all, models.
Errors in converting to ESM-SnowMIP format because of the above.	N/A	Some models. Results analysis team.
Hard-coded errors		
Bug in model use of site longitude	Unrealistically low albedo with consequences on snow mass and melt (AT)	JULES-GL7, JULES-UKESM

Bug in transmission of SW radiation through canopy	Investigated slow melting behaviour of model after evaluation data became available (MT)	SURFEX-ISBA
Model SWE limited to a maximum of 1000 mm	SWE limited to 1000 mm (AT)	MATSIRO
Unintentional decoupling of snow surface and atmosphere	Snow did not melt at Weissfluhjoch in some summers (AT)	HTESSEL-ML
Bug in partitioning of SW radiation into direct and diffuse	Unrealistically high albedo values (AT)	Crocus
Bug in the output of liquid water content	Found unrealistically small liquid water content values when compared ESM-SnowMIP results with other simulations (MT)	HTESSEL, EC-EARTH
Inconsistent use of snow area fraction when calculating snow depth and SWE	Snow density varied instead of being fixed (AT)	MATSIRO
Many variations in output file formats; wrong variable name; variations in the interpretation of the ESM-SnowMIP definition of output variables; different sign conventions.	N/A	One or more in adjacent list for most, if not all, models.
Errors in converting to ESM-SnowMIP format because of the above.	N/A	Some models. Results analysis team.

Figures

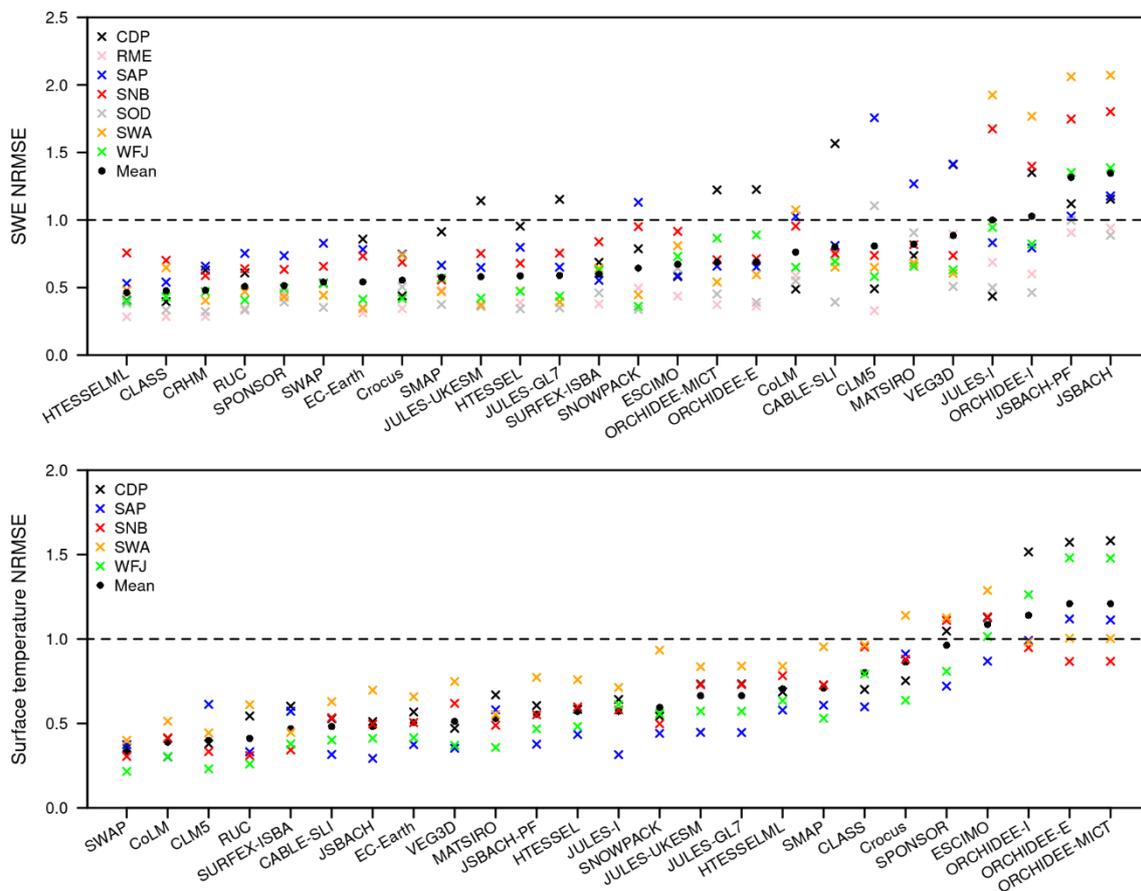


Figure 1: Model ranking by normalised root mean square errors of snow water equivalent and surface temperature. The site names are shortened as follows: CDP = Col de Porte, SAP = Sapporo, RME = Reynolds Mountain East, SNB = Senator Beck, SOD = Sodankylä, SWA = Swamp Angel and WFJ = Weissfluhjoch.

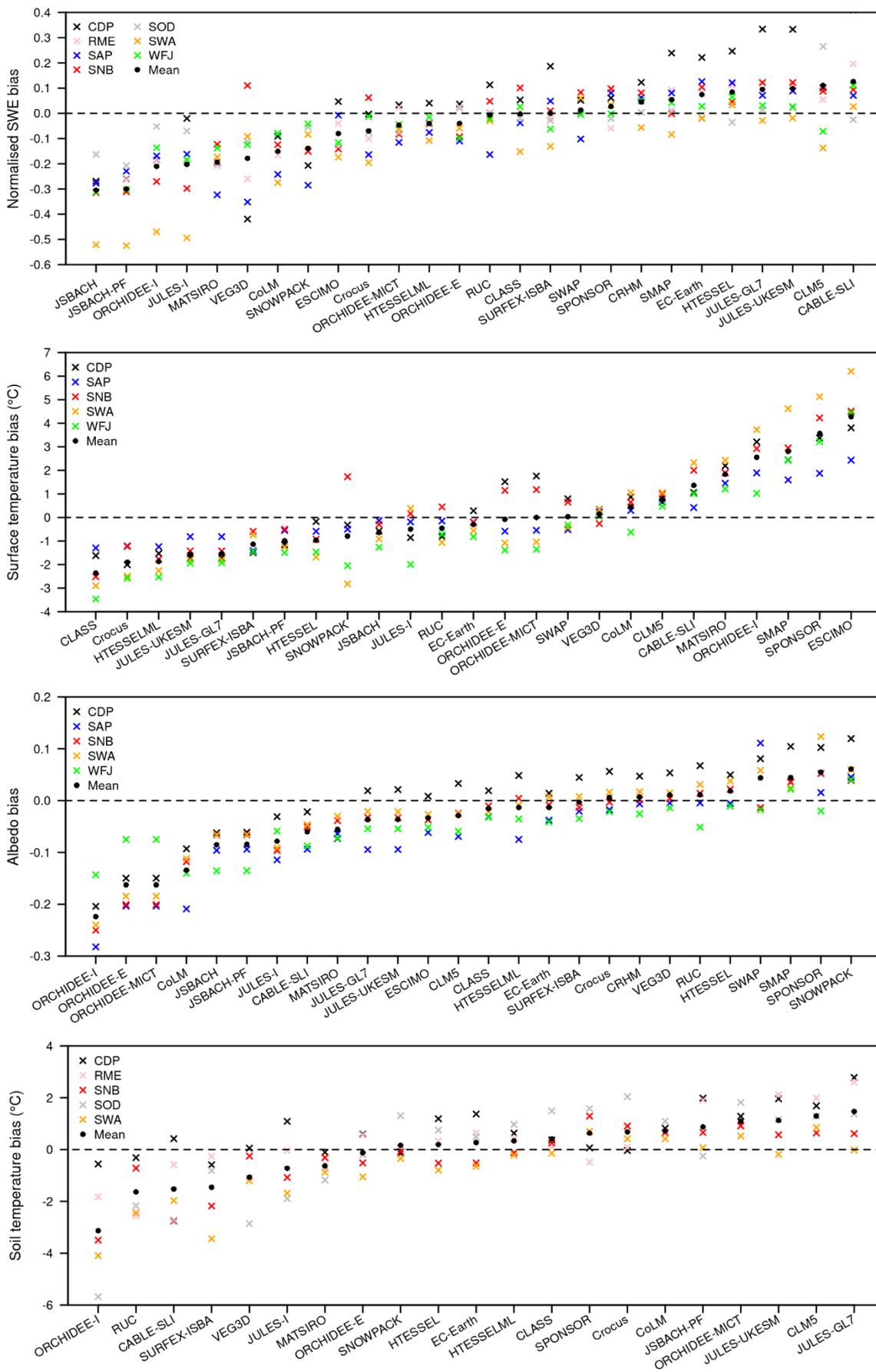


Figure 2: Model ranking by biases from negative to positive. Following the prevalent

convention, negative biases denote model underestimates. SWE biases are normalised by measured mean yearly maxima. JSBACH soil temperature cold biases (ranging from -6°C to -12°C and averaging -9°C) are outside the range of the plot. The site names are shortened as in Figure 1.

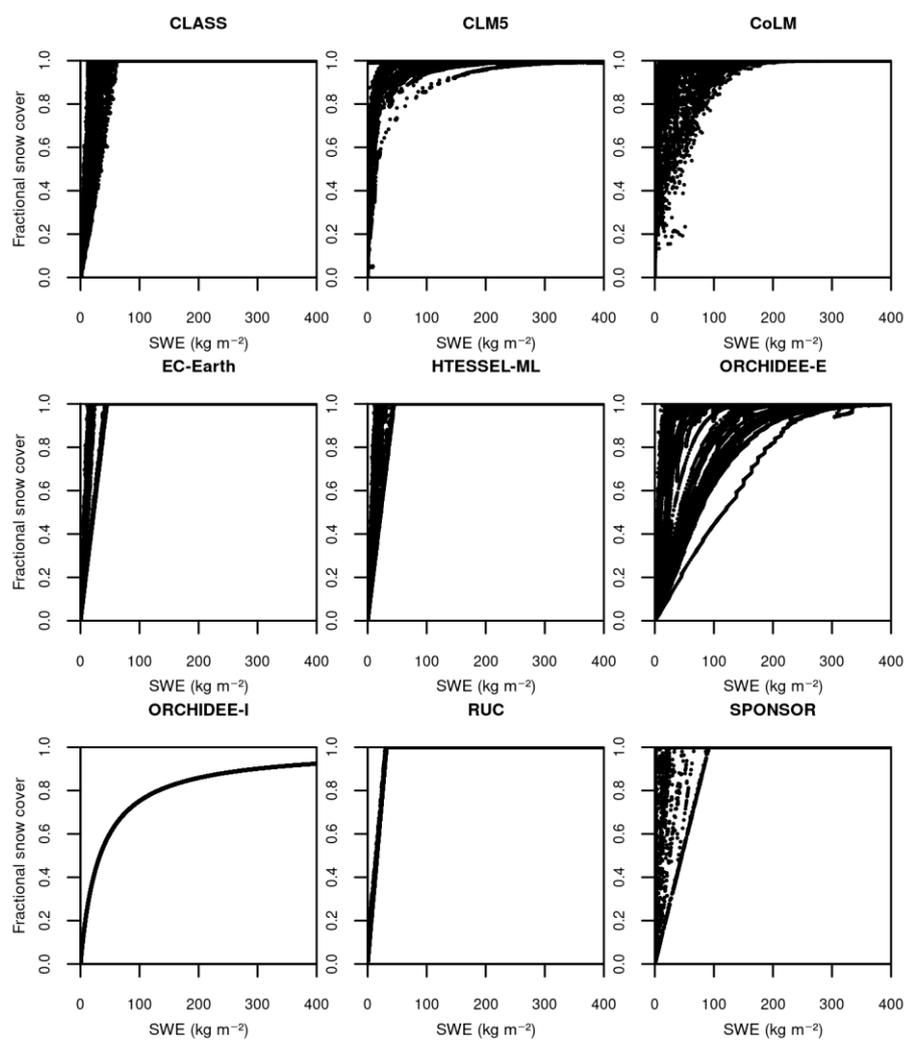


Figure 3: Fractional snow cover (SCF) as a function of SWE at Col de Porte for models that did not switch off their sub-grid parametrizations or impose complete snow cover. HTESSEL is not shown as it is the same as HTESSEL-ML. ORCHIDEE-MICT did not force SCF = 1, but values were missing from the file provided for evaluation.

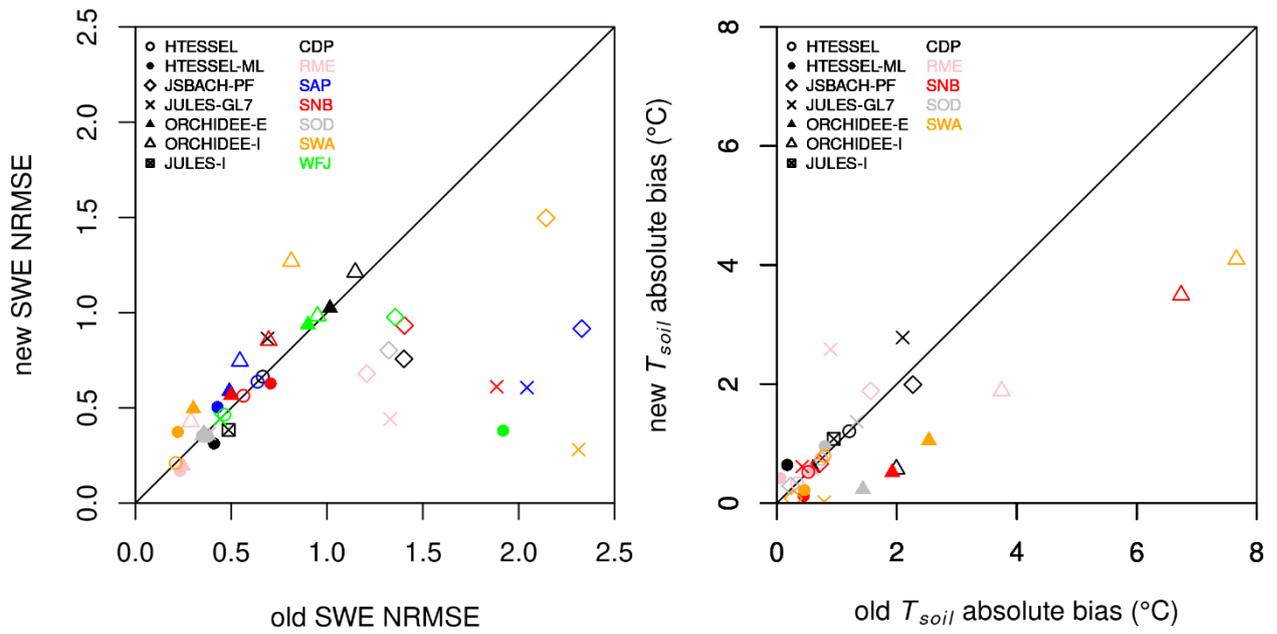


Figure 4: SWE NRMSE and soil temperature (T_{soil}) absolute bias before and after resubmission for selected models. The site names are shortened as in Figure 1.