

Referee #1

General comments: The authors of the manuscript (ms.) have tested six non-calibrated snow models at one mountain location by varying the time-resolution and origin of model forcing. The quality of meteorological forcing is indeed an important element in snow modelling, and the authors have here examined the sensitivity of snow model performance to varying input data quality. The manuscript is quite well-written, and the illustrations are mostly clear and understandable. The research idea is based on straightforward model testing, i.e. running an ensemble of models with different forcing data and calculating statistics on how the models' performance (evaluated against point snow observations) vary. However, the two main weaknesses of the ms. are in my opinion that:

- 1) the model evaluation is made at only one single site (Torgnon)
- 2) there is no discussion or treatment of uncertainty in solid precipitation measurements, which directly affect snow amounts in the comparison.

Although the ms. may provide an interesting case study for those interested in the specific models and the local site, the ms. does not have, in my opinion, large enough impact and interest for a wider audience to warrant publication in HESS. After reading the ms. twice, I felt in the end that I did not get much relevant information out of it for my modelling work, in another place, another country. As the authors themselves state (p. 22, lines 8-9): "This study offers some hints on this research topic".

We reply to the two points highlighted by the reviewer in the text below, as these comments have been reported and expanded by the reviewer in the "Specific comments" section (see below). In particular, regarding the issue of the interest of the paper for a wider audience not specifically working on the models and site considered in this study, we address this in detail in our replies (mainly 1. and 5.) to better clarify the aims of our work, the results that can be exported to other snow models and sites, and the still-open issues, left for future investigation.

Specific comments:

1. As snow is often spatially very inhomogeneously distributed, normally the utility in snow modelling is to get a grasp of this spatial variability. Thus, simulating snow in just one point has limited relevance, mostly restricted to snow process studies. In another point, the authors' results and model ranks might be changed significantly. As the authors themselves state, on p. 2, line 15: "Snow models are generally evaluated at a number of sites"; on p. 26, lines 2-3: "Further analysis at other test sites would be useful to explore the extent to which our results could be generalized to different situations or models".

We see the reviewer's point and we are aware that the choice of a single study-site can be highlighted as a limitation, however we have some motivations to support this choice.

The strengths of our work are, in our opinion, the analysis of a multi-model ensemble, representative of different degrees of complexity of snow models, and the analysis of a wide range of possible meteorological forcing datasets, to explore in detail the response of the models to forcings with different characteristics and resolutions in time and space. When planning this large, collaborative experiment, we carefully considered the choice of the site where to perform our analysis. The site we finally selected is quite unique as it provides high quality data in particular for precipitation (in most cases poorly measured in high elevation sites) and it is affected by low wind speeds, so that the snow-drift effect is limited. The combination of these two conditions is rare in high-elevation mountain

measurement sites, nevertheless it is essential if we want to reduce the uncertainties on the input data. Repeating this effort in multiple test sites, for example in other alpine sites at different elevations and latitudes, or at non-alpine sites, i.e. in the Arctic) would certainly expand the results provided by the paper but at the cost of larger uncertainties in the forcings which propagate across the modeling exercise and complicate the interpretation of the model outputs. Reducing the uncertainty on the “control” forcings is a prerequisite, in our case, to better separate the error due to model structure from the error due to the forcing. In this context, the selected site has represented for us the most appropriate benchmark for the aim of the paper. Extending the investigation to other test sites with less “optimal” forcing would be of great interest but, in our opinion, it should be addressed in a separate paper.

We hope to have clarified the motivations underlying our choice. We believe that this study, shedding light on the impacts of the model complexity and of the accuracy of the forcing on snow simulations, could be of interest for the readers of Hydrology and Earth System Sciences involved in catchment hydrology, snow modelling and snow and water resources management.

- 2. The authors note on p.4 line 1-12: “the uncertainty on snow simulations due to the forcing can be comparable to or even larger than the uncertainty”. They also refer on p.7 line 2 to Kochendorfer et al. (2017), who assess and provide algorithms to deal with the undercatch of solid precipitation. However, no effort is made to discuss, assess or correct the precipitation measurements at Torgnon station for the undercatch and/or examine the sensitivity of the authors’ results for the inherent uncertainties in the observation-based model precipitation input (their CTL experiment with “optimal forcing”).**

Following the reviewer’s suggestion we analyzed in more detail the uncertainty associated with the observed precipitation and in particular the undercatch of snow which is common in mountain areas. The primary cause for snow precipitation undercatch is related to wind speed, with the amount of precipitation measured by a precipitation gauge relative to the actual amount of precipitation decreasing with increasing wind speed.

We quantified the wind-induced precipitation measurements errors by applying the method described in Kochendorfer et al. (2017). This method, derived by comparing precipitation measurements from unshielded and shielded (reference) gauges, consists in calculating a catch efficiency (CE), function of *air temperature* and *wind speed*, so that its inverse (CE^{-1}) can be used to correct actual precipitation data. The method has been specifically developed for OTT Pluvio2 gauges, i.e. the same type as that used at the Torgnon site.

Figure 1 shows the cumulated total precipitation at the Torgnon site measured by the precipitation gauge (black) compared to the precipitation adjusted with the Kochendorfer method (blue).

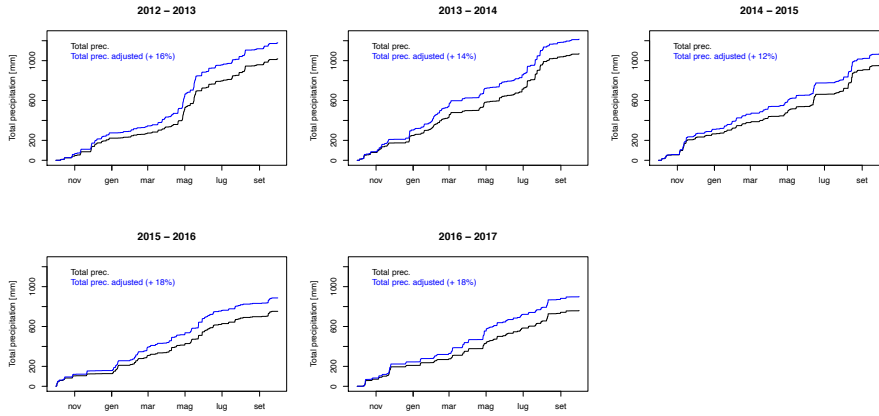


Figure 1 Cumulated total precipitation at the Torgnon site measured by the OTT Pluvio2 precipitation gauge (black) compared to the precipitation adjusted with the Kochendorfer method (blue).

The adjusted cumulated total precipitation exceeds the measured precipitation by 16% in average over the 5 snow seasons.

As the correction of total precipitation directly affects the amount of solid precipitation, we tested the effects of such correction on snow model simulations. We performed an additional experiment (CTL_prc-adj) in which the model forcing is the same as in the CTL run except for total precipitation, which is now *adjusted*, and snowfall which is now calculated from the *adjusted* total precipitation.

Figure 2 shows the results for the SNOWPACK model, and it displays the simulated snow depth (upper panel) and snow water equivalent (bottom panel) obtained in the CTL and in the CTL_prc-adj runs compared to observations.

In all snow seasons the snow depth and the snow water equivalent are remarkably overestimated in the CTL_prc-adj experiment compared to both observations and the CTL run. The additional snowfall input derived from the precipitation adjustment leads to an excess of snow accumulation on the ground which can be quantified in an average snow depth bias of 0.17 m compared to the -0.001 m bias in the CTL run. The RMSE is double in the CTL_prc-adj run compared to the CTL run (Table 1).

Given that the precipitation adjustment method itself is affected by its own uncertainties, and given that the application of the precipitation adjustment leads to a worsening in the snow model performances, we prefer to employ the original precipitation measurements as forcing in the snow model experiments. The discussion of the uncertainty of precipitation measurements and the effect of the precipitation adjustment on snow simulations has been included in the Appendix of the revised manuscript and summarized in the main text.

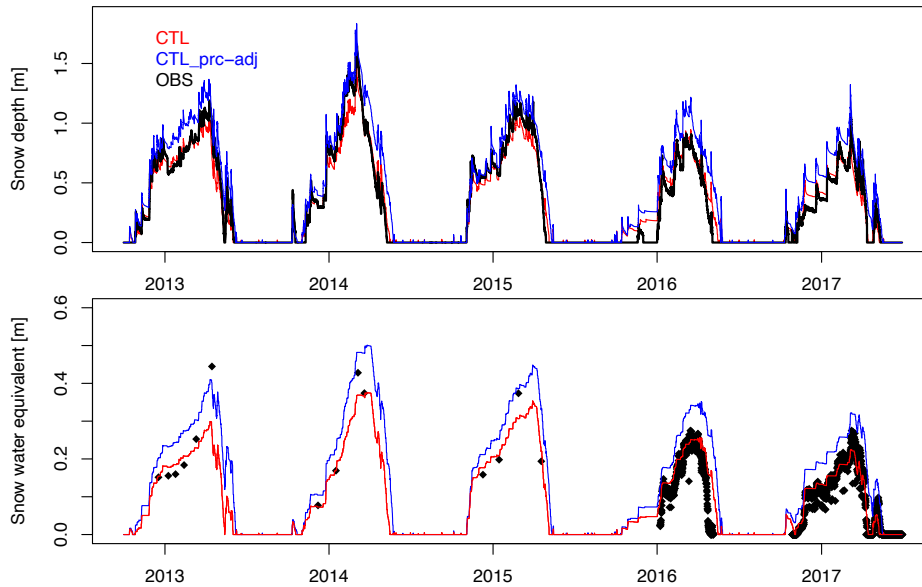


Figure 2 Snow depth (upper panel) and snow water equivalent (lower panel) simulated by the SNOWPACK model when the adjusted total precipitation forcing is employed (CTL_prc-adj) compared to the control run (CTL) and observations.

Table 1. SNOWPACK model RMSE and bias for the simulated snow depth and snow water equivalent variables in the control run (CTL) and in the CTL_prc-adj experiment.

	Snow depth		SWE	
	RMSE [m]	BIAS [m]	RMSE [m]	BIAS [m]
CTL	0.10	-0.001	0.04	0.02
CTL_prc-adj	0.20	0.170	0.10	0.09

3. The authors claim that a bias adjustment of forcing data leads to more precise results (p.9. lines 29-31). This seems to me like a rather trivial point.

The sentence in the manuscript reads “The last two experiments [...] investigate if it is possible to improve the performances of snow models [...] by applying two simple bias-correction methods to adjust air temperature and hence the amount of solid precipitation with respect to the total one.”

The idea here is to check if:

- Correcting temperatures only (and keeping all the other variables unchanged except for solid and liquid precipitation whose partition depends on temperature) can improve the model snow simulations
- Very simple bias correction methods (such as the lapse rate correction and the subtraction of the mean bias) can be sufficient to improve model performances or more sophisticated techniques are necessary.

We rephrased the sentence in the text to clarify the meaning.

4. The linear interpolation of shortwave radiation e.g. in the TIME-12h case, causing the large deviations of +97 W/m² (p. 18, lines 3-4), is an unrealistically simplistic way to make the interpolation. In real modelling practice, I suppose most of us would use a sinus-curve form or something like that. Consequently, the issue here is more of poor modeling practice than lower time-resolution. This is only mentioned in section

Discussions (p.24, line 10), but would have been best to put into practice already in the authors' study.

Following your suggestion we tested a more realistic way of estimating the 30 minute incoming shortwave radiation when only the measurements at 00:00 and 12:00 are available, i.e. as in the TIME-12h experiment. We employed the potential (clear-sky) incoming shortwave radiation (Knauer et al., 2018) at 30 minute temporal resolution and at the coordinates corresponding to the Torgnon station, and the surface station SWIN measurements at 12:00.

For each day of the year, the 48 daily values of potential radiation are rescaled according to the observed SWIN value at 12:00, to obtain an “estimated SWIN” (see Figure 3).

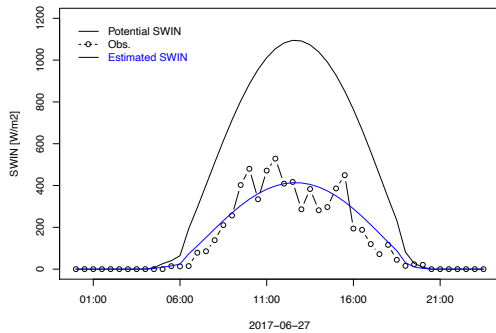


Figure 3. Measured shortwave incoming radiation (SWIN) at the Torgnon site for the day 26 June 2017 (points), potential SWIN for that day (solid black line), “estimated SWIN” from the scaling of the potential SWIN on the value registered at h 12:00.

The advantage of this method compared to the linear interpolation method is that the difference between the estimated and the observed SWIN radiation averaged over the full period is almost cancelled out, from +97 W/m² when using the linear interpolation method to -0.87 W/m² when using the method based on the scaling of the potential radiation.

In light of this result we run a new experiment TIME-12h-SWIN-POT, in which the forcing is the same as the one employed in the original TIME-12h experiment except for the shortwave incoming radiation, which is now obtained with the potential radiation method.

Figure 4 shows the results of the TIME-12h-SWIN-POT experiment compared to that of the original TIME-12h experiment, the CTL run and observations, for the SNOWPACK model and for the snow depth variable. The use of the potential radiation remarkably improves the agreement with observations, reducing the RMSE with respect to observations to a value which is comparable to the CTL run (Table 2).

The results of the TIME-12h-SWIN-POT experiment have been reported in the revised version of the manuscript and the effects of the two different interpolation methods (one based on the linear interpolation of the measurements and the other based on the scaling of the potential radiation) on the snow simulations have been discussed in the main text and more extensively in Appendix C.

SNOWPACK model

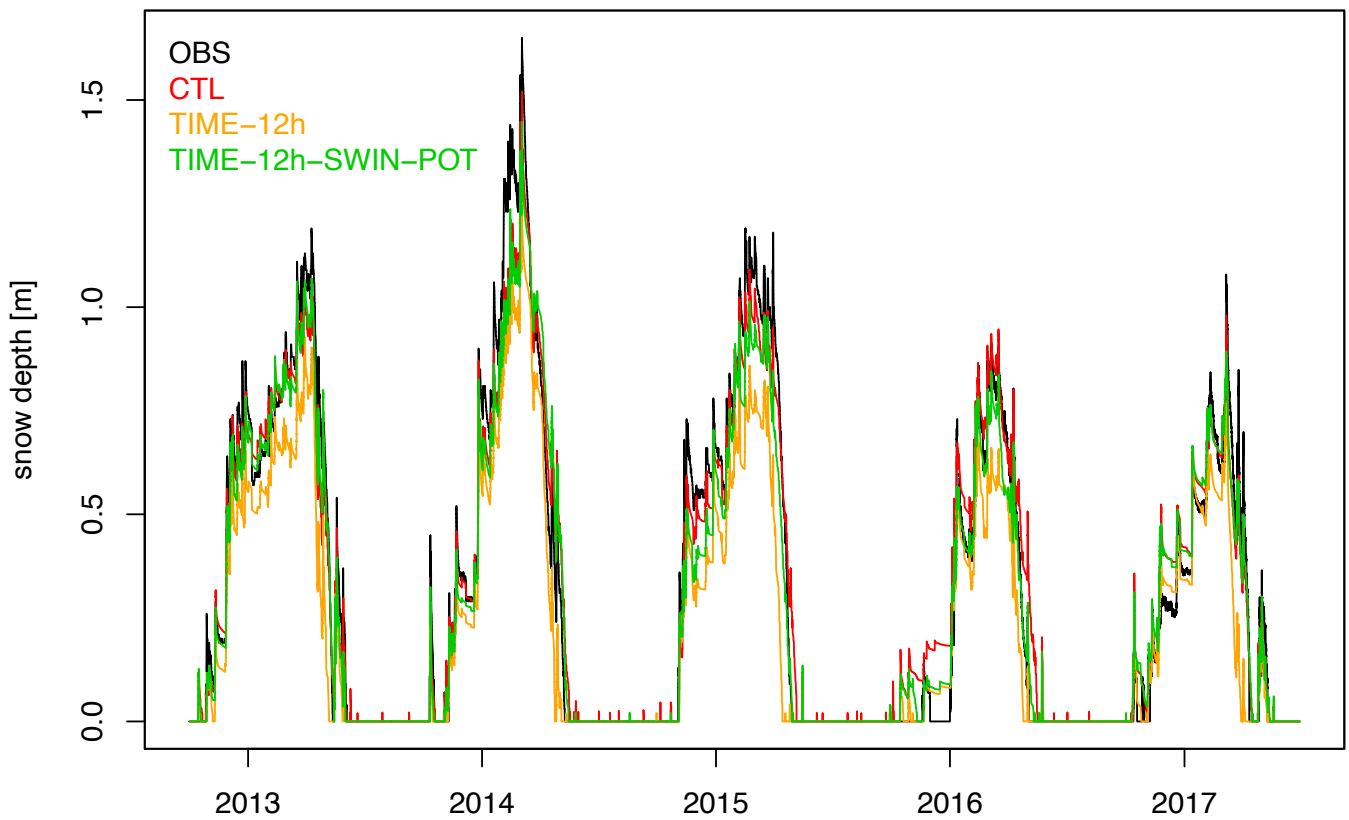


Figure 4 Snow depth simulations obtained with the SNOWPACK model for the experiment TIME-12h-SWIN-POT compared to TIME-12h, the CTL run and observations.

Table 1 SNOWPACK model RMSE, BIAS and Pearson Correlation for the simulated snow depth in the CTL run, the TIME-12h and the TIME-12h-SWIN-POT experiments, compared to observations.

	Snow depth		
	RMSE [m]	BIAS [m]	Pearson correlation
CTL	0.10	-0.001	0.97
TIME-12h	0.21	-0.016	0.93
TIME-12h-SWIN-POT	0.11	-0.036	0.97

5. The point-specific biases and errors of the MeteoIO, GLDAS, ERA5 and ERA Interim, described in Sections 5.4 and 5.5 for the single Torgnon site, emphasize the weakness of this case study: the model evaluation results are difficult to generalize outside this site, where things and biases could be very different. Also, compensating errors may occur in the models which improve model performance. In other words, one gets “right results for wrong reasons” (p.13, lines 14-16; p. 21, lines 4-5; also p.24, line 17).

We disagree that the biases and errors highlighted for MeteoIO, GLDAS, ERA5, and ERA-Interim at the Torgnon site emphasize the weaknesses of the work for two reasons. Concerning the presence

of biases, the aim of the paper is indeed to test how snow models respond to inputs which might be affected by large uncertainties and errors. Concerning the reviewer's remark on the difficulty of generalizing the results outside the area of study we addressed this point by testing the reanalysis products against observations over the Greater Alpine Region (GAR), to observe the spatial distribution of the temperature and precipitation biases and see if they are consistent at the mountain range scale.

ERA5, ERA-Interim and GLDAS temperatures have been averaged over the months October-June and over the years 1980-2014 (except for GLDAS which is available since 2000 only, so the averages have been calculated over the period 2000-2014), and then compared to the observational dataset EOBS (version 13, Haylock et al., 2008). EOBS is a daily gridded data set at 0.25° resolution, based on the European Climate Assessment and Data set station measurements.

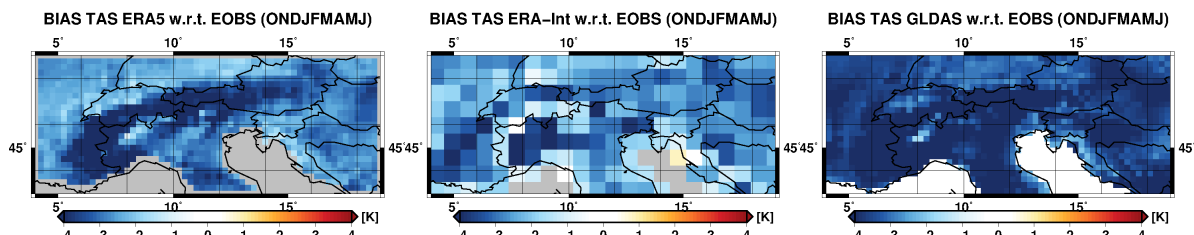


Figure 5 BIAS of ERA5, ERA-Interim and GLDAS air temperatures with respect to EOBS observations over the Greater Alpine Region. Temperatures have been averaged over the months from October to June and over the period 1980-2014 in the case of ERA5 and ERA-Interim, over the period 2000-2014 in the case of GLDAS. .

ERA5 and GLDAS temperature biases are large and negative over the entire GAR (Figure 5). GLDAS bias is especially strong and it exceeds -4°C in most of the region. ERA5 bias is larger at high elevation than in lowlands. Compared to ERA5 and GLDAS, ERA-Interim temperature is in better agreement with observations, with mainly negative bias across the region and values close to zero (both positive and negative values) except at the mountain ridges in Western Alps.

Regarding precipitation, it is well known that standard surface station gauges have problems in capturing snowfall and thus they underestimate total precipitation in mountain areas. Similarly, also observational-based dataset such as EOBS have been found to suffer the underestimation of precipitation at high elevations (Turco et al., 2013). To overcome this problem, instead of using observation-based datasets as a reference, we evaluate precipitation differences with respect to a reanalysis, which inherently takes into account orographic effects. Figure 6 shows the ERA5 and GLDAS October-to-June accumulated precipitation differences relative to ERA-Interim (ratio) over the periods 1980-2014 and 2000-2014 respectively (GLDAS is available since 2000). Also in this case ERA5 spatial pattern is homogeneous over the Alpine range, with ERA5 showing consistently more precipitation than ERA-Interim in the mountain areas. Concerning GLDAS, we need to clarify that, while working on this response to reviewers, we noticed an error in the method which we used to perform the temporal interpolation of the original data and to derive the 30-minute resolution precipitation forcing for the “GLDAS” experiment. The error has now been fixed and the snow model runs driven by the GLDAS reanalysis have been repeated with the correct forcing. The updated results have been reported in the revised version of the manuscript. With this correction, moreover, we found a much better agreement of the GLDAS precipitation with both Torgnon station measurements and with the ERA-Interim reanalysis over the Greater Alpine Region (Figure 6, right panel).

Overall, this analysis providing information on the spatial variability of the temperature and precipitation biases in the reanalysis products over the Alpine region broadens the perspective beyond the specific case of the Torgnon site. The biases at the Torgnon site result generally coherent with those found at the mountain range scale, although the magnitude of the bias can vary across the region

and with the elevation. This analysis addresses the question of how the bias in the main forcing variables (temperature and precipitation) at the Torgnon site can be generalized at larger scale, and in particular over the Alpine region.

The main results of this analysis are reported in the “Discussion” section of the revised manuscript, while the plots are reported in the Appendix D.

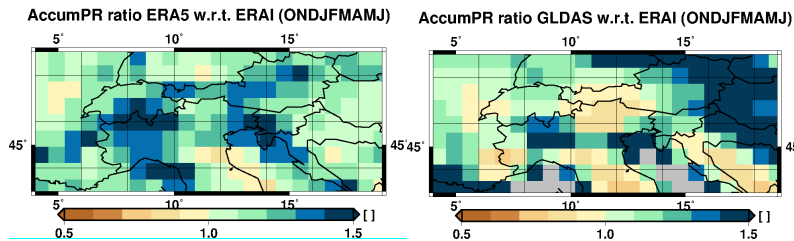


Figure 6 ERA5 and GLDAS relative differences with respect to ERA-Interim for the October-June accumulated precipitation over the periods 1980-2014 and 2000-2014 respectively.

6. Normally, the authors' results would be compared to previous studies in the section “Discussion”. However, this three-page section only has one (!) single reference to other published studies. It also repeats many things already pointed out previously in the ms.

We thank the Reviewer for this comment, the Discussion has been extensively modified avoiding repetitions and including the comparison of our results with those obtained in similar studies previously published in literature.

Other points:

7. p.7 lines 22-23: provide values for the vertical gradients used here.

The information has been added in the text, thank you

8. Fig. 2. The meaning of dashed line circles in the Taylor graph should be explained.

Guidance on how to interpret Taylor diagrams has been added in the text, thank you.

9. The Appendix is short (a little figure and table, and five lines of text) and could be easily added to main text in Section 3.2.

In the revised version of the manuscript the Appendix has been expanded. It now includes i) the discussion on the uncertainty of the total precipitation measurements at the Torgnon station; ii) information on how the Meteo-IO interpolated data have been derived; iii) the discussion of the impact of the time interpolation method for SWIN in the TIME-12h experiment; iv) the discussion of the biases of the reanalyses over the Greater Alpine Region.

10. The ms. contains a lot of numbers in tables 3-4, and displaying them more readerfriendly, like in the nice Fig. 6, would be good.

Thanks for the suggestion, however in this case we preferred to keep the numerical values of data in the classic table form.

Although limited to a single site and short on mechanistic explanations, the evaluations of several models in simulating snow mass, depth and density with several forcing datasets in this paper are of value (another with similar aims that should be cited is <https://journals.ametsoc.org/doi/full/10.1175/JHM-D-15-0013.1>).

We thank the reviewer for this suggestion, the citation has been included in the revised paper.

Measurements of outgoing shortwave and longwave radiation are mentioned but not used in the model evaluations; these might provide more insight.

We agree that the evaluation of modelled outgoing shortwave and longwave radiation could provide interesting additional insights, as suggested by the reviewer. These variables are not provided by all the models considered in this study: they are missing completely in the simplest model considered, S3M, while the SMASH model provides the outgoing longwave radiation only. Therefore, we evaluated the simulated outgoing shortwave and longwave radiation at the surface for all remaining models (SNOWPACK, GEOTOP, HTESSEL and UTOPIA).

Figure 7 shows the difference (left) and the scatterplot (right) of the simulated and observed daily-averaged outgoing shortwave radiation for the SNOWPACK, GEOTOP, HTESSEL and UTOPIA models in the CTL run.

All the models tend to moderately underestimate the outgoing shortwave radiation, with SNOWPACK and HTESSEL showing the best agreement with observations both in terms of bias and of coefficient of determination R^2 , compared to GEOTOP and UTOPIA. Differences between the simulated and observed outgoing shortwave radiation are mainly dependent on the representation of the albedo. These results suggest to check, in the UTOPIA model, the albedo, which is function of surface temperature and snow age as well as in the GEOTOP model, which is function of snow age and grainsize.

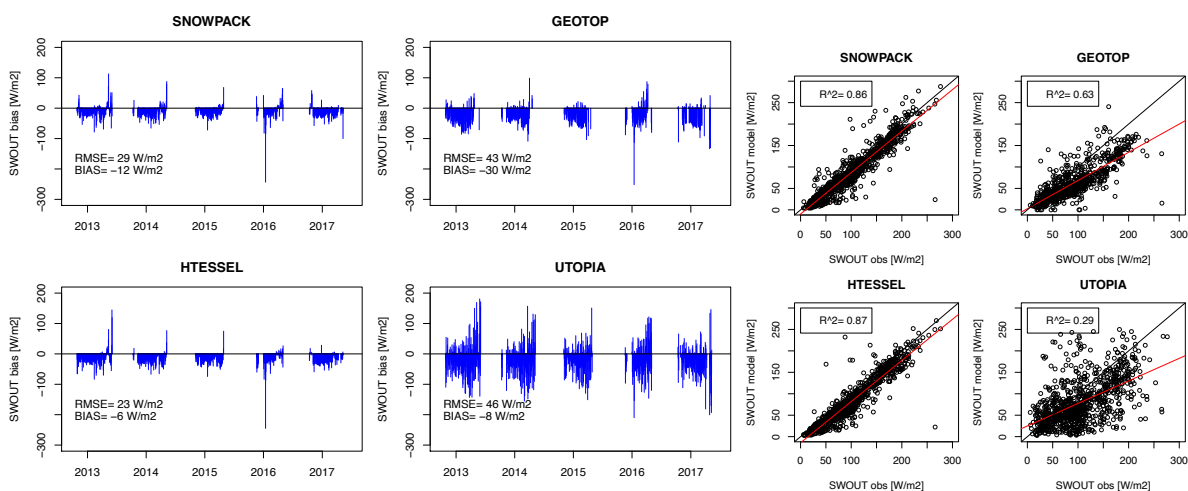


Figure 7 Difference (left figure) and scatterplot (right figure) of the simulated and observed outgoing shortwave radiation (CTL experiment) at the Torgnon site. Only four out of the six models considered in the paper provide the outgoing shortwave radiation among the output variables.

Similarly to Figure 7, Figure 8 shows the difference (left) and the scatterplot (right) of the simulated and observed daily-averaged outgoing longwave radiation for the SNOWPACK, GEOTOP, HTESSEL, UTOPIA and SMASH models in the CTL run. The simulation of the net longwave radiation mainly affects the representation of the snow-melt dynamics. Both SNOWPACK and SMASH underestimate the outgoing longwave radiation, causing an excess in the snowpack energy

available for the melting. However, none of the two models remarkably underestimates the snow depth, so other mechanisms might compensate for this behavior. GEOTOP, HTESSEL and UTOPIA outgoing longwave radiation do not show systematic biases.

Since the outgoing shortwave and longwave radiation outputs are not provided by all the six models considered in the study and we lack the information for simpler models, we preferred not to report this analysis in the manuscript. We thank the reviewer for the suggestion.

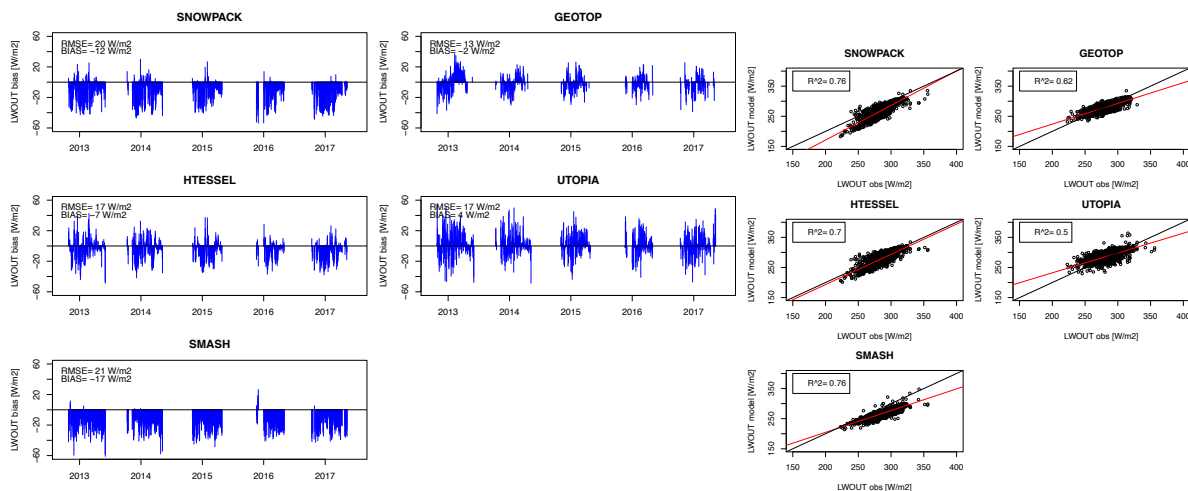


Figure 8 Difference (left figure) and scatterplot (right figure) of the modelled and observed outgoing longwave radiation (CTL experiment) at the Torgnon site. Only five out of the six models considered in the paper provide the outgoing longwave radiation among the output variables.

- **p2, line 28 There is no need for feedbacks for early differences in snow simulations to persist throughout the winter if the imposed conditions remain too cold for melt.**

We thank the reviewer for the suggestion; the text has been modified accordingly.

- **p3, line 22 delete “air” in “open air sites”**

Done, thank you.

- **p3, line 25 Specifically, Rutter et al. (2009) found that benefits from calibration at forest sites did not transfer to nearby non-forested sites. Direct calibration at the non-forested sites would almost certainly have improved simulations.**

We thank the reviewer for the suggestion; the text has been modified accordingly.

- **p6, line 18 I do not think that wind direction is needed to force the snow models, and please clarify whether any of them use surface temperature.**

Indeed none of the models employs the wind direction, thanks for the correction. With “surface temperature” we actually meant “ground temperature at 2 cm depth”, and this variable is needed by the SNOWPACK model only. We have now clarified both points in the text.

- **p7, line 1 “both liquid and solid fractions” means that total precipitation is measured, not separate snowfall and rainfall.**

Yes, exactly. At the Torgnon station the total precipitation amount is measured. We modified the text to better clarify it, thank you.

- **p7, line 23 After reading this, I expected Appendix A to give details of the vertical gradients used for temperature and precipitation interpolation.**

Yes, we added this information in the text, thanks for the suggestion.

- **p8, Table 1 Information on the elevation of the reanalysis grid points would be interesting here or in the text. Also, how much gap-filling was required in the station data?**

Information on the elevation of the reanalyses at the Torgnon gridpoint, as well as the % of missing value for each input variable provided by the Torgnon station have been added in Table 2.

- **p9 I am confused by SWIN-CLS. If R is measured radiation and SWIN is modelled clearsky radiation, I don't see where the MSG cloud masks are being used. If R is incident solar radiation in cloudy conditions, isn't equation 2 the wrong way round?**

Yes, the equation was wrong and it has been corrected, many thanks for pointing it out. MSG cloud mask is used to identify the radiometers under clouds and compute an average attenuation factor. We have better explained this in the text.

- **p9, 13 Linear interpolation of sampled radiation fluxes rather than solar elevation-based interpolation of accumulated fluxes will be biased. How do average fluxes compare? (briefly mentioned in 5.3 and turns out to be a source of error)**

We agree that linear interpolation of sampled shortwave radiation fluxes introduces errors in the forcing data, which can be large when the sampling time step is 12h (TIME-12h experiment). For this case, we tested a different method to estimate 30 minutes shortwave incoming radiation (SWIN) from 12h samplings, based on the rescaling of the potential radiation with respect to the measurement at 12:00.

The details of this exercise, as well as the comparison between the linearly interpolated SWIN and the modified SWIN forcing, are provided in the reply to Referee#1, at point 4. In summary, while the linearly interpolated SWIN forcing shows an average bias of +97 W/m² compared to observations in the period of investigation, the modified SWIN has significantly reduced the bias to a value close to zero (-0.87 W/m²), so the average flux is conserved. We ran an additional experiment (TIME-12h-SWIN-POT) using the modified SWIN forcing (see details in the reply to Referee#1, point 4).

A detailed discussion of this point and the results of the new experiment have been included in the manuscript in Sections 4 and 5.3

- **p10 Is the partitioning of total precipitation into snowfall and rainfall only applied to station measurements (in which case I expected to read about it in section 3) or also to the reanalyses, even though they provide separate snowfall and rainfall?**

We applied the same method to separate rainfall and snowfall for all the forcing datasets, including reanalyses. We better clarified it in Section 3 and 4.

- **p12, line 31 Even a model that could account for impurities would not do so in this case because dust deposition was not provided as an input.**

Yes, thank you for this comment. We modified the text as suggested

- **p18, line 33 MeteIO errors are relatively small for temperature and snowfall, but errors in other forcing variables are not shown. Correct spelling of "systematically" throughout**

The MeteIO forcing biases with respect to the Torgnon measurements are relatively small on average not only for temperature and snowfall but for all variables (Figure 9). In order to explain the discrepancies between the simulated SWE and snow depth in the MeteIO experiment and observations (see i.e. Table 4 in the manuscript), we investigated the temporal variability of the air temperature bias (Figure 10) and we related it to the simulated SWE and snow depth in the MeteIO experiment (Figure 11). The temperature bias is about -1°C on average over the considered time

period, however in winter the cold bias is generally stronger and it can reach values exceeding -4°C (Figure 10). Concerning MeteIO-driven snow model simulations, the main issue is the overestimation of snow depth in winter (in selected snow seasons) and in spring (in all seasons). A plausible explanation for these errors is that colder-than-observed winter temperatures might favor the development of a cold snowpack which melts too slowly. Consequently, the models tend to overestimate the snow at surface and to predict a delayed ablation date.

We added these comments in Section 5.4 of the manuscript.

The typo “systematically” has been corrected, thank you.

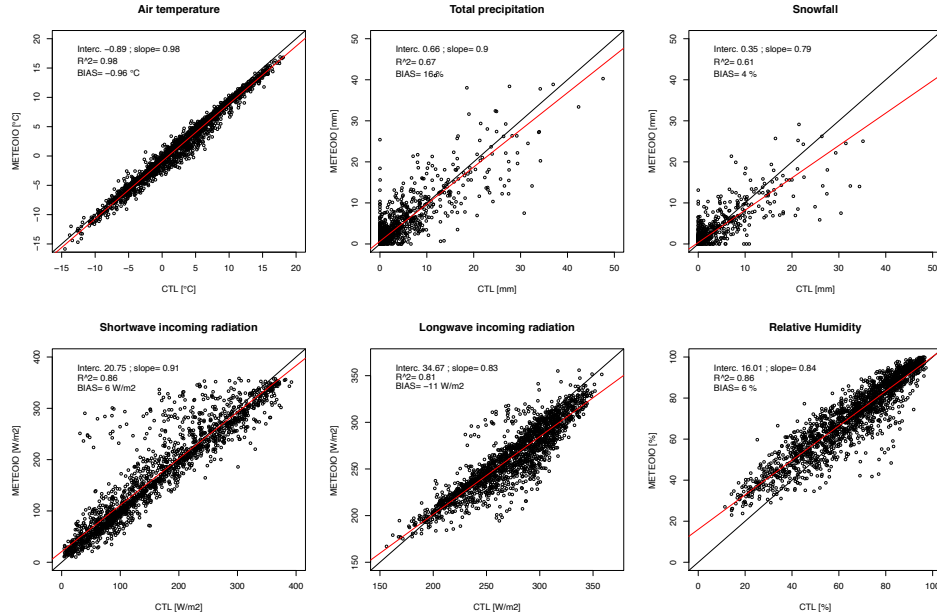


Figure 9. Scatterplot of the meteorological forcing of the MeteIO experiment with respect to the CTL run.

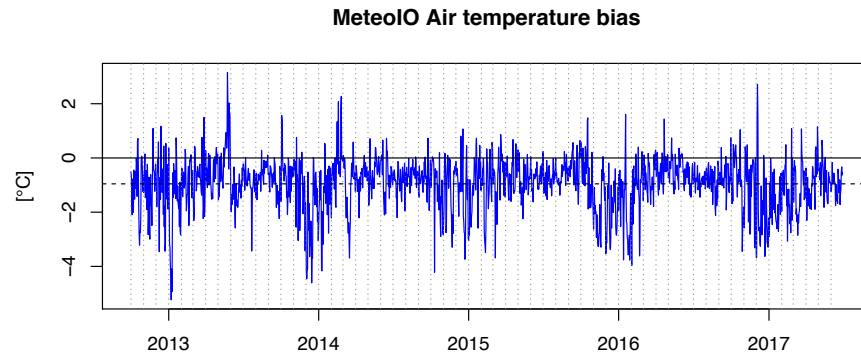


Figure 10 MeteIO air temperature bias with respect to Torgnon station measurements.

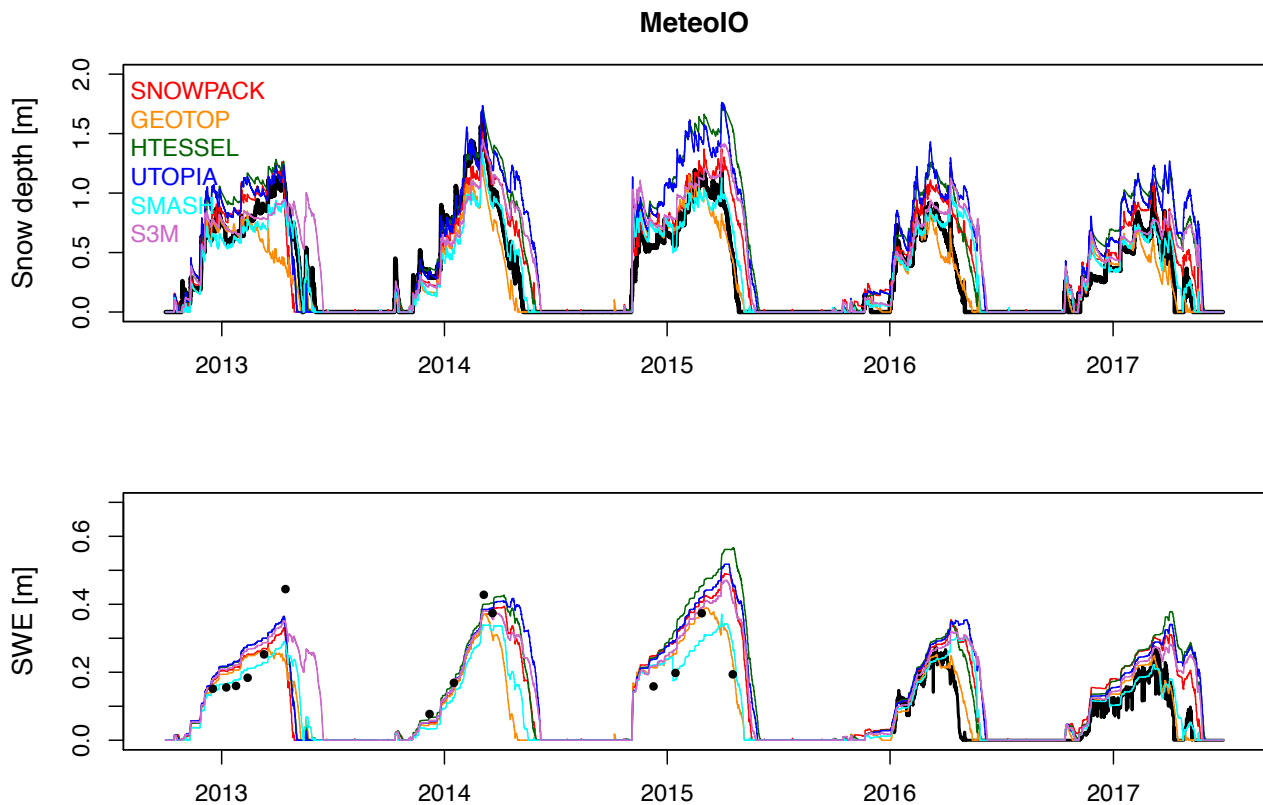


Figure 11 Simulated snow depth (top) and snow water equivalent for the MeteolO experiment compared to observations.

References

- Haylock, M., Hofstra, N., Klein Tank, A., Klok, E., Jones, P., and New, M.: A European daily high-resolution gridded data set of surface temperature and precipitation for 1950–2006, *J. Geophys. Res.-Atmos.*, 113, D20119, <https://doi.org/10.1029/2008JD010201>, 2008
- Knauer, Jürgen, Tarek S. El-Madany, Sönke Zaehle, and Mirco Migliavacca. "Bigleaf—An R package for the calculation of physical and physiological ecosystem properties from eddy covariance data." *PloS one* 13, no. 8 (2018).
- Kochendorfer, J., Nitu, R., Wolff, M., Mekis, E., Rasmussen, R., Baker, B., Earle, M. E., Reverdin, A., Wong, K., Smith, C. D., Yang, D., Roulet, Y.-A., Buisan, S., Laine, T., Lee, G., Aceituno, J. L. C., Alastrué, J., Isaksen, K., Meyers, T., Brækkan, R., Landolt, S., Jachcik, A., and Poikonen, A.: Analysis of single-Alter-shielded and unshielded measurements of mixed and solid precipitation from WMO-SPICE, *Hydrol. Earth Syst. Sci.*, 21, 3525–3542, <https://doi.org/10.5194/hess-21-3525-2017>, 2017.
- Turco, M., Zollo, A. L., Ronchi, C., De Luigi, C., & Mercogliano, P. (2013). Assessing gridded observations for daily precipitation extremes in the Alps with a focus on northwest Italy. *Natural Hazards & Earth System Sciences*, 13(6).

Sensitivity of snow models to the accuracy of meteorological forcings in mountain environment

Silvia Terzago¹, Valentina Andreoli², Gabriele Arduini³, Gianpaolo Balsamo³, Lorenzo Campo⁴, Claudio Cassardo², Edoardo Cremonese⁵, Daniele Dolia⁴, Simone Gabellani⁴, Jost von Hardenberg^{6,1}, Umberto Morra di Cella⁵, Elisa Palazzi¹, Gaia Piazzzi^{4,7}, Paolo Pogliotti⁵, and Antonello Provenzale⁸

¹Institute of Atmospheric Sciences and Climate, National Research Council, Torino, Italy

²Department of Physics and Natrisk center, University of Torino, Italy

³European Centre for Medium-Range Weather Forecasts, Reading, UK

⁴CIMA Research Foundation - International Centre on Environmental Monitoring, Savona, Italy

⁵Environmental Protection Agency of Aosta Valley, Aosta, Italy

⁶Department of Environment, Land and Infrastructure Engineering, Politecnico di Torino, Italy

⁷IRSTEA, Hydrology Research Group, UR HYCAR, 92761 Antony, France

⁸Institute of geosciences and earth resources, National Research Council, Pisa, Italy

Correspondence: Silvia Terzago (s.terzago@isac.cnr.it)

Abstract. Snow models are usually evaluated at sites providing high-quality meteorological data, so that the uncertainty in the meteorological input data can be neglected when assessing the model performances. However, high-quality input data are rarely available in mountain areas and, in practical applications, the meteorological forcing used to drive snow models is typically derived from spatial interpolation of the available in-situ data or from reanalyses, whose accuracy can be considerably lower.

5 In order to fully characterize the performances of a snow model, the model sensitivity to errors in the input data should be quantified.

In this study we test the ability of six snow models to reproduce snow water equivalent, snow density and snow depth when they are forced by meteorological input data with gradually lower accuracy. The SNOWPACK, GEOTOP, HTESSEL, UTOPIA, SMASH and S3M snow models are forced, first, with high-quality measurements performed at the experimental site of Torgnon, located at 2160 m a.s.l. in the Italian Alps (control run). Then, the models are forced by data at gradually lower temporal and/or spatial resolutions, obtained i) by resolution, obtained by i) sampling the original Torgnon 30-minute time series at 3, 6, and 12 hours, ii) by spatially interpolating neighboring in-situ station measurements and iii) by extracting information from GLDAS, ERA5, ERA-Interim reanalyses at the gridpoint closest to the Torgnon station site. Since the selected models are characterized by different degrees of complexity, from highly sophisticated multi-layer snow models to simple, empirical, single-layer snow schemes, we also discuss the results of these experiments in relation to the model complexity.

15 ~~Results show that~~ The results show that, when forced by accurate 30-min resolution weather station data, the single-layer, intermediate-complexity snow models HTESSEL and UTOPIA provide similar skills as the more sophisticated multi-layer model SNOWPACK, and these three models show better agreement with observations and more robust performances over different seasons compared to the lower complexity models SMASH and S3M. All models forced by 3-hourly data provide similar skills as the control run while with the use of 6- and 12-hourly temporal resolution forcings we generally observe may

lead to a reduction in model performances ~~except for the SMASH model which~~ in case the incoming shortwave radiation is not properly represented. The SMASH model generally shows low sensitivity to the temporal degradation of the input data. Spatially interpolated data from neighboring stations and reanalyses result to be adequate forcings, provided that temperature and precipitation variables are not affected by large biases over the considered period. ~~A~~ However, a simple bias-adjustment technique applied to ERA-Interim temperatures ~~however~~, allowed all models to achieve similar performances as ~~in~~ the control run. ~~All models irrespectively~~ Irrespectively of their complexity, all models show weaknesses in the representation of the snow density.

Copyright statement. TEXT

1 Introduction

A wide range of snow models with different degrees of complexity have been developed for hydrological applications, avalanche risk forecasting ~~or~~ and climate studies. Some of them are also integrated within modelling chains for numerical weather forecasts or climate modelling. The degree of complexity of the snow schemes depends on the specific purpose for which they have been developed (Magnusson et al., 2015). Simple temperature-index snow models are employed in applications requiring a coarse estimate of snow depth or snow water equivalent. Physical, energy-balance, but still rather simple snow models are often used in complex modelling chains, i.e. in numerical weather prediction systems and in Earth System models, to limit the computational costs. Sophisticated multi-layer snow models are typically used to reconstruct the vertical structure of the snowpack with a high level of detail and high accuracy, as needed for avalanche warning applications.

Snow models are generally evaluated at a number of sites providing high-quality forcing and verification data. Extensive literature documents the underlying physics and the performances of single snow models (e.g. Dutra et al., 2010; Vionnet et al., 2012; Bartelt and Lehning, 2002), and several studies compare a limited number of snow models ~~to~~ with each other (Boone and Etchevers, 2001; Kumar et al., 2013). A few large intercomparison studies benchmarked multiple snow models, including the PILPS2d, PILPS2e, Rhone-Agg, SNOWMIP and SNOWMIP2 coordinated intercomparison projects.

PILPS2d (Slater et al., 2001; Schlosser et al., 2000) and PILPS2e (Bowling et al., 2003) aimed at evaluating snow water equivalent (SWE) simulations provided by different land surface schemes (LSS) in a Russian and a Swedish snow-dominated catchment, respectively. PILPS2d evaluated twentyone land surface schemes forced by 18 years of observed meteorological data from a grassland catchment in Russia, ~~with the aim~~ to investigate the reasons for models scatter in the output snowpack variables. Weaknesses in reproducing mid-season ablation were shown to produce systematic scatter among between the models. Albedo and fractional snow cover were both key variables for an accurate representation of the amount of energy absorbed by the snowpack. ~~Indeed, the~~ The ablation during the early snow season is another major source of divergence among between models: in early winter a thin snow cover is highly sensitive to changes in the forcings and the resulting differences in snow-

pack conditions tend to persist throughout the whole snow season ~~owing to internal feedback processes~~ if temperatures remain too cold for melt.

PILPS2e showed the difficulty of reproducing spring melting. Errors in winter snow sublimation mainly impacted the runoff simulations, while the retention of meltwater within the snowpack affected the timing of the peak in runoff rather than its magnitude. For both PILPS2d and PILPS2e the differences in model complexity did not fully explain the differences in model results.

The Rhône-AGG experiment (Boone et al., 2004) employed 15 LSSs to address the impact of the model structure and of the spatial resolution of the forcing data on the simulations of the water balance. LSSs with an explicit (bulk or multi-layer) snow scheme provided better SWE simulations than LSSs with a composite snow scheme (i.e with a mixed snow-soil-vegetation layer). LSSs with composite snow schemes showed ~~too~~ early snow ablation and early run-off peaks compared to observations, owing to missing representation of key processes such as ripening and ~~owing to~~ inadequate representation of albedo and thermal conductivity in a mixed snow-soil/vegetation layer. SWE was strongly affected by the spatial resolution of the meteorological forcing. In fact, when high-resolution meteorological forcings were aggregated from 8 km to a coarser grid of 1° (about 69 km) the simulated SWE was reduced by 25-60% in 13 out of 15 LSSs. A single model explicitly considering subgrid elevation effects on the forcing was found to minimize the impact of scaling on the simulated snow water equivalent.

SnowMIP (Etchevers et al., 2002, 2004) performed an intercomparison ~~among of~~ snow models of different complexity, used for different applications, including hydrology, global circulation models, snow monitoring, snow physics, avalanche forecasting, with the aim of identifying key processes for each application. Model complexity was found to have a strong impact on the simulation of the net longwave radiation, which strongly affects snow melt dynamics. Models relying on the explicit simulation of the internal snow processes ~~better~~ represented snow surface temperature and the longwave radiation budget more accurately. On the contrary, model complexity had smaller impact on the net shortwave radiation, whose accuracy was dependent on the simulation of albedo. Complex models taking into account snow microstructure were able to properly represent the albedo variability (as a function of grain size and type), but also simple snow models with an appropriate parameterization ~~for of~~ albedo dynamics guaranteed reliable estimates of this variable.

SnowMIP2 (Rutter et al., 2009) built upon SnowMIP and focused on the simulation of snowpack properties in forested areas compared to open ~~air~~ sites, across different climatic conditions. Single models showed low correlations between different years in forested sites, and low correlations also between forested and open sites, suggesting that no single best model for all years and all sites could be easily identified. Calibration ~~can~~ allowed to reduce root mean square error (RMSE) in forested sites ~~but similar improvements did not apply to~~ but the benefits from calibration at forested sites did not transfer to nearby non-forested sites.

The mentioned studies shed light on the critical snow processes that produce the largest differences between LSS simulations. However they could not clearly define an optimal set of parameterizations for a given application, ~~as for example such as~~ numerical weather predictions and climate simulations, or the minimum level of model complexity needed to achieve satisfactory skills in a given application (Slater et al., 2001). A step forward in this direction ~~has been made~~ was obtained by employing a single model with several options to represent each of the most snow-relevant processes, and then testing the effect of

parameterizations with different degrees of complexity on the skill of the model (Essery et al., 2013; Clark et al., 2011). ~~Best~~ The best results were obtained with models ~~with~~ having a prognostic representation of snow albedo and density, with at least a simple representation of water retention and refreezing in the snowpack. ~~Major outcomes on the identification of the key processes to be represented in global climate models to improve snow simulations at the large scale are expected from the ongoing coordinated modelling~~ The ongoing coordinated initiative ESM-SnowMIP (Krinner et al., 2018) is expected to provide important information on the key snow processes that should be included in Global Climate Models.

A common characteristic among past model intercomparison initiatives is the interest in testing the skills of the models in experimental sites where high-quality meteorological forcings are available, ~~in order~~ to perform a controlled evaluation of the ~~models~~ model performances. However, such context does not represent the typical conditions occurring in practical applications, where snow models are run over large climate model grid cells, and they are coupled to atmospheric models that likely provide biased driving data (Essery et al., 2013). Moreover, reliable modelling of snowpack dynamics in mountain regions is hindered by the high spatial and temporal variability of the meteorological forcings, entailing that observations and reanalysis data at a given location are scarcely representative of the conditions of the surrounding area. A recent review paper on the European mountain cryosphere (Beniston et al., 2018) states that disentangling the uncertainties related to the model structure from those related to the meteorological input data is one of the major challenges for the snow modelling at the catchment scale, ~~namely the scale~~ relevant for hydrological applications. A sensitivity analysis performed on a single, physically-based snow model showed that the uncertainty on snow simulations due to the forcing can be comparable to or even larger than the uncertainty due to the model structure (Raleigh et al., 2015). ~~The~~ That analysis also showed that biases in the forcing data have a larger effect than random errors, ~~with the magnitude of the biases resulting more important than their probability distribution~~. Building on the results of previous studies we now ~~consider~~ expand the perspective by considering an ensemble of snow models with different degrees of complexity and we investigate their sensitivity to the quality of the meteorological forcing, with the aim of providing ~~practical~~ information on their performances when they are forced with inputs at gradually lower temporal and/or spatial resolution.

We devised a set of experiments with six ~~different~~ snow models with different degrees of complexity in the Alpine measurement site of Torgnon, located at 2160 m a.s.l. in Aosta Valley, Italy. ~~This site has been selected because it provides high-quality meteorological measurements together with a characterization of snowpack properties in terms of depth, mass and surface temperature~~. First, we evaluate each model forced by accurate station measurements at 30-minute temporal resolution (we refer to this as "optimal" forcing). Second, we ~~evaluate each model~~ test the response of each model when forced by data at gradually lower temporal resolution and/or lower accuracy, ~~by employing~~. To this end, we employ data from spatial interpolation of neighboring station measurements and from three gridded global reanalyses, and ~~extracting~~ we extract the meteorological time series at the gridpoint closest to the Torgnon station. The ~~set of multi-model, multi-forcing simulations finally allows to discuss the relationship between model complexity~~ site of Torgnon has been selected because it provides high-quality meteorological measurements, in particular for precipitation which is usually poorly measured in high elevation sites, and a detailed characterization of snowpack properties in terms of depth, mass and surface temperature. Moreover the Torgnon site usually experiences low wind speeds, so that the snow-drift effect is very limited. The combination of these three conditions is

rare in high-elevation measurement sites but essential to reduce the uncertainties on input and validation data, and to allow for a reliable estimation of the error due to model structure. Repeating this effort at multiple test sites, for example in other alpine sites at different elevations and latitudes, or at non-alpine sites (i.e. in the Arctic) would expand the results provided by the present paper. Of course, this would come at the cost of larger uncertainties in the forcings, which propagate across the modeling exercise and complicate the interpretation and the comparison of the model outputs. For this reason, we leave this more complex investigation for a separate paper. Here we employ a multi-model and ~~model-sensitivity-to-errors~~ multi-forcing framework to i) assess the performances of each snow model when forced with inputs at gradually lower temporal and ~~inconsistencies in the meteorological forcings~~ or spatial resolution, ii) discuss the relation between model performances and model complexity, and iii) provide model users with information for practical applications.

This paper is structured as follows: Sect. 2 presents the snow models employed in the study while Sect. 3 describes the station of Torgnon and the datasets employed for the experiments. Section 4 describes in detail the set of 12 devised experiments and, for each experiment, the method employed to derive the forcing. Section 5 focuses on the evaluation of snow model outputs against observations, finally Sections 6 and 7 discuss the results and draw the conclusions.

2 Snow models

The six models considered in this study, together with a compact overview of their characteristics, are listed in Table 1 and summarized in the following.

SNOWPACK is a highly sophisticated, multi-purpose snow and land-surface model, with a detailed description of the mass and energy exchange between the snow, the atmosphere and optionally ~~with~~ the vegetation cover and the soil. It provides a detailed description of snow properties including weak layer characterization (Stoessel et al., 2009), phase changes and water transport in snow (Hirashima et al., 2010). A particular feature is the treatment of soil and snow as a continuum with a choice of a few up to several hundred layers (Bartelt and Lehning, 2002).

GEOTOP 2.0 is a sophisticated, snow and hydrological process-based model. Its strength is an integrated approach that takes into account the interactions between hydrological, cryospheric and geomorphological processes (Endrizzi et al., 2014). The snowpack evolution is dynamically managed by the model through a snow layering scheme which splits and merges the layers depending on their mass. ~~Moreover, the~~ The model takes into account also snow metamorphism and water percolation into the snowpack.

HTESSEL is the land-surface model of the European Centre for Medium-Range Weather Forecasts (ECMWF) Integrated Forecasting System (IFS), controlling the evolution of the snow and soil fields and the exchanges of heat and moisture between the land surface and the atmosphere above (Balsamo et al., 2009). HTESSEL includes a process-based single-layer snow scheme to represent the grid cell fraction (tile) that is covered by ~~the~~ snow (Dutra et al., 2010). In this scheme, the snowpack is characterized by a prognostic temperature, mass, density and albedo, updated at each time step. The liquid water content is diagnosed based on the other snow fields (temperature, density and mass), allowing to represent the interception of rainfall by the snowpack and internal melting/refreezing processes (Dutra et al., 2012).

Table 1. Features of snow models in terms of model complexity following Slater et al. (2001), snow albedo (α) parameterization, explicit representation of meltwater retention and refreezing in the snowpack (M_w) and a main reference.

Snow model	Complexity	α^*	M_w	Reference
SNOWPACK	multi-layer	111	Yes	Bartelt and Lehning (2002)
GEOTOP	multi-layer	011	Yes	Endrizzi et al. (2014)
HTESSEL	single-layer	110	Yes	Dutra et al. (2012)
UTOPIA	single-layer	110	Yes	Cassardo (2015)
SMASH	up to 3 layers	110	No	Piazzì et al. (2018, 2019)
S3M	single layer	010	No	Boni et al. (2010)

*The three digits 1 and 0 represent the dependence or not of the albedo parameterization respectively on surface temperature, snow age, grainsize. 000 means fixed albedo.

UTOPIA is a land-surface process model representing the physical processes at the interface between surface, vegetation and soil layers, including a scheme which accounts for the main processes occurring in the snowpack (Cassardo, 2015). The snowpack is considered as a single homogenous layer placed upon land surface and its mass, thermal and hydrological balances are analyzed. The model takes into account the partition of soil coverage fractions (bare soil, vegetated soil, soil or vegetation covered by snow) and is able to simulate snow water equivalent, depth, density, albedo and coverage. Snow metamorphism is parameterized.

SMASH is a two-layer snow model that reproduces some of the main physical processes occurring within the snowpack, including accumulation, density dynamics, melting, sublimation, radiative balance, heat and mass ~~exchanges~~exchange (Piazzì et al., 2019). The model can be coupled with multivariable data assimilation schemes (Piazzì et al., 2018, 2019) allowing the joint assimilation of several snow-related observations to produce SWE and runoff estimates. ~~With the aim of facilitating the~~
To facilitate the implementation of the assimilation algorithms, the complexity of the modelling scheme is ~~accordingly~~-limited (e.g., liquid water storage and refreezing process are neglected). ~~It is noteworthy that in~~In the present study no assimilation scheme has been implemented in SMASH (open-loop configuration).

S3M is a spatially distributed, empirical snow model requiring only few input variables (precipitation, temperature, incoming shortwave radiation and air humidity) to compute the water mass conservation equation and to produce a first estimate of SWE (Boni et al., 2010). A second, optional, independent estimate of the SWE field, obtained combining spatial interpolation of surface snow depth observations and MODIS snow cover, is assimilated into the snow model using a nudging scheme. The result of the data assimilation is an updated SWE map exploiting different sources of information, modeling, remote sensing and surface stations network measurements. In the use of the model for the experiments proposed in this paper the assimilation scheme is switched off and the model runs in ~~open-loop~~open-loop configuration.

In the proposed experiment all the models are used in their default configurations, so no special tuning of the model parameters is ~~done~~made to improve the results over the Torgnon site. All models calculate snow water equivalent and snow density as primary variables, while snow depth is derived from them.

3 Study site and data

3.1 Torgnon station data

Meteorological forcing data are provided by a fully-equipped weather observation station located at Torgnon, 2160 m a.s.l. (45°50' N, 7°34'E) in the Aosta Valley, Western Italian Alps. The experimental site belongs to the ICOS (IT-Tor, www.icos-ri.eu/) and LTER (lter_eu_it_077, www.data.ltereurope.net/deims/site) networks and it is described in detail by Galvagno et al. (2013), Filippa et al. (2015) and Piazzzi et al. (2019). The location is a subalpine grassland, an abandoned pasture located a few kilometres from the village of Torgnon. The site is characterized by an intra-alpine semi-continental climate, with mean annual temperature and precipitation of 3.1°C and 880 mm respectively (Galvagno et al., 2013). During the cold season most of precipitation falls as snow and, on average, from the end of October to late May, the site is snow covered with snow depths reaching 90-120 cm (Galvagno et al., 2013). Wind-induced phenomena are limited in this site, since it experiences low winds, with an average half-hourly wind speed of 1.6 ± 1.3 m/s over the 2012-2014 period.

The station measures all the input variables needed to force the snow models, including air temperature, total precipitation, shortwave (SWIN) and longwave (LWIN) incoming radiation, wind ~~direction and~~ speed, relative humidity, surface pressure, ~~surface temperature and ground temperature at 2 cm depth (the last variable is employed by the SNOWPACK model only).~~ These variables are measured at high frequency, and then aggregated at 30-minute temporal resolution. Precipitation measurements are performed with an OTT Pluvio2 Weighing Rain Gauge, which employs a weight-based technique to measure both liquid and solid fractions ~~-(i.e the total precipitation amount).~~ This is a consolidated technique that provides higher confidence on the reliability of precipitation data than standard rain gauges ~~(?)~~. ~~(Kochendorfer et al., 2017a). Despite the station being equipped with a reliable pluviometer and it is exposed to low wind speeds, possible issues of precipitation undercatch can be present. The uncertainty associated with precipitation measurement has been estimated and the impact of the uncertainty of the precipitation input on snow model simulations has been assessed and discussed in Appendix A.~~ As the OTT pluviometer has been operational since mid-2012, in our analysis we consider the dataset spanning the period from October 1st, 2012 to June 30th, 2017, covering five complete snow seasons.

The Torgnon station provides also snow-related variables useful for model evaluation, including snow depth measurements, obtained by an ultrasonic distance sensor, surface temperature, snow and soil temperatures at different depths, outgoing short-wave and longwave radiation, all of them available at 30-minute resolution. Snow density and snow water equivalent are measured manually in snow pits several times per snow season during dedicated field campaigns. During the analysis period 20 manual measurements of snow density and snow water equivalent are available. Additionally, since January 2016 snow water equivalent is automatically monitored by a Campbell CS725 sensor, that passively measures the attenuation of naturally existing electromagnetic radiation (Potassium-40 and Thallium-208) emitted from the soil or bedrock below the sensor. The higher the water content of the snow pack, the higher the attenuation of the radiation. The measure is performed every 6 hours and averages the SWE over an area of about 100 m². Combining automatic snow water equivalent measurements and the corresponding snow depth measurements, additional daily snow density estimates useful for model validation have been derived for the last two snow seasons.

3.2 Spatial interpolation of meteorological forcings from neighboring stations

The spatial interpolation of ground meteorological observations represents one of the most commonly used practices in the operational applications of hydrological models. In order to test the performances of the models in this condition, an interpolated dataset has been generated for the Torgnon monitoring site by using the MeteIO library (Bavay and Egger, 2014). Meteorological data from six neighboring stations have been interpolated over a squared digital elevation model of 16 km² with a grid resolution of 50 ~~meters-m~~ centered on the coordinates of Torgnon ([Appendix B](#) Fig. B1 and Tab. [A+B1](#)). The algorithm used for the interpolation is the inverse distance weight (IDW) as first choice for all the meteorological variables. The interpolation accounts also for vertical gradients of both temperature and precipitation, assuming constant lapse rates of -6.5°C/km for air temperature and +8.5 mm/km for precipitation. Further details are provided in [Appendix AB](#).

3.3 Reanalysis data

In many remote mountain areas, in-situ observations to force snow models are unavailable. In this study we explore the use of reanalysis datasets extracted at the Torgnon gridoint.

GLDAS (Global Land Data Assimilation System) is a global dataset exploiting satellite and ground-based observational data combined with advanced modelling and data assimilation techniques in order to generate optimal fields of surface variables (Rodell et al., 2004). In particular, the GLDAS-2.1 archive used in this study contains 36 land surface fields from January 2000 ~~to present time and updated regularly~~ at 0.25° (lon/lat) spatial and 3-hour temporal resolutions (Rui and Beaudoin, 2018).

ERA-Interim (Dee et al., 2011) is a global reanalysis including a variety of 3-hourly surface parameters describing atmospheric and land-surface conditions, and 6-hourly upper-air parameters covering the troposphere and stratosphere. ERA-Interim has spatial resolution of 0.75°, at the latitude of Torgnon corresponding to about 59 km in the zonal and 83 km in the meridional direction. This coarse grid, which is comparable to those of state-of-the-art global climate models, implies a smooth representation of the topography and coarse information on climate variables.

ERA5 (Hersbach and Dee, 2016) is the latest ECMWF global reanalysis product, providing data at higher resolution than ERA-Interim, both in space (30 km) and in time (1 hour). ERA-5 uses one of the most recent versions of the Earth system model and data assimilation methods applied at ECMWF and modern parameterizations of Earth processes compared to older versions used in ERA-Interim. With respect to ERA-Interim, ERA5 has also an improved global hydrological and mass balance, reduced biases in precipitation, and refinements of the variability and trends of surface air temperature (Hersbach and Dee, 2016).

4 Experimental design

We devised a set of twelve experiments at the Torgnon site employing snow models in stand-alone mode, i.e. in which the meteorological forcing is prescribed. The list of experiments is summarized in Table 2. The first experiment is a control run (CTL) in which the models are forced by optimal input data provided by the Torgnon station at 30-minute temporal resolution.

Table 2. Overview of the experiments and their characteristics in term of forcing data, temporal and spatial resolutions and gap-filling data employed where necessary. [For reanalysis datasets, the elevation of the gridpoint closest to the Torgnon station is reported.](#)

Experiment	Forcing	Temporal resolution	Spatial Resolution	Gapfilling
CTL	Torgnon station (2160 m a.s.l.)	30'	point	ERA-Interim ERA-Interim*
RAD-ERA-Interim	CTL except SWIN and LWIN from ERA-Interim in case of snowfall	30'	point	ERA-Interim
SWIN-CLS	CTL except SWIN from Clearsky algorithm	30'	point	ERA-Interim
TIME-3h	Torgnon station	3h	point	ERA-Interim
TIME-6h	Torgnon station	6h	point	ERA-Interim
TIME-12h	Torgnon station	12h	point	ERA-Interim
MeteoIO	Six stations close to Torgnon (see Appendix A)	1h	point	none
GLDAS	GLDAS-2.1 (2297 m a.s.l.)	3h	25 km	none
ERA5	ERA5 (2302 m a.s.l.)	1h	30 km	none
ERA-Interim	ERA-Interim (1480 m a.s.l.)	3h	80 km	none
ERA-Interim-LR	ERA-Interim, lapse-rate correction of air temp. temperature	3h	80 km	none
ERA-Interim-BIAS	ERA-Interim, bias adjustment of air temp. temperature	3h	80 km	none

* % of missing values for each variable: air temperature 0.24%; surface air pressure 1%; wind speed 1.65%; total precipitation 0.25%; shortwave incoming radiation 0.33%; relative humidity 0.24%; longwave incoming radiation 0.31%.

This run allows to test the accuracy of the models in describing the temporal evolution of the snow-related variables in optimal conditions, namely when high-quality, high-frequency point measurements are available.

Experiments ~~LWIN-ERA-Interim~~ RAD-ERA-Interim and SWIN-CLS assess the sensitivity of the models to the radiation input. As most stations, [the](#) Torgnon site is equipped with an unheated radiation sensor, which is likely to provide unreliable measurements when getting obstructed by snow during snowfall events. Therefore, in the experiment RAD-ERA-Interim we take into account the shading of the radiation sensor in case of snowfall by replacing radiometer measurements with ERA-Interim reanalysis data. In the third experiment, SWIN-CLS, we employ external SWIN data resulting from the clear sky radiation (Yang et al., 2001, 2006) attenuated through the cloud masks from the Meteosat Second Generation ([MSG](#)) satellite in the following way. For each of the 34 radiometers in the Aosta Valley an averaged attenuation factor F is computed as:

$$F = \frac{1}{N} \sum_{i=1}^N \frac{R_{st}^i}{SWIN^i} \quad (1)$$

where N is the number of cloud-covered stations [determined from the MSG cloud mask](#), R_{st}^i is the measured radiation at the i^{th} station and $SWIN^i$ is the corresponding modeled radiation in ~~clear-sky~~ clear-sky condition. The incident solar radiation in cloudy conditions at the location j is given by:

$$\underline{R^j} = SWIN^j_{inc} = R^j F \quad (2)$$

Experiments TIME-3h, TIME-6h and TIME-12h investigate the sensitivity of the models to the temporal resolution of the meteorological forcing, since the temporal resolution of many available datasets is coarser than that employed in the CTL run. We ~~basically employ~~ employ the Torgnon data every 3, 6 and 12 hours since October 1st, 2012 time 00:00 UTC +01:00, and linearly interpolate them at the 30 minute time step for all variables except for total precipitation. Precipitation is accumulated over 3, 6 or 12 hour time periods and the totals are equally distributed among the corresponding 30 minutes subperiods. Incoming shortwave radiation is linearly interpolated at the 30-minute time step for all experiments, i.e. TIME-3h, TIME-6h and TIME-12h. However, when we apply linear interpolation to derive the forcing for the TIME-12h experiment we obtain poor SWIN estimates, with a large difference between the estimated and the CTL average SWIN flux (+97 W/m²). In order to better estimate the SWIN forcing for the TIME-12h experiment we employ a method based on the potential (clear-sky) radiation at 30 minute temporal resolution (Knauer et al., 2018) at the site of Torgnon. For each day of the year, the 48 values of potential radiation are rescaled according to the observed SWIN at 12:00. With this method the estimated average SWIN flux is comparable to that of the CTL forcing, with a difference of -0.87 W/m^2 , showing a remarkable improvement with respect to the use of the linear method (more details provided in Appendix C). We run the TIME-12h experiment twice, either employing the SWIN derived from the linear interpolation method (TIME-12h-LIN) or that derived from the potential radiation method (TIME-12h-SWINPOT)

Four additional experiments, namely MeteoIO, GLDAS, ERA5 and ERAI test the case in which no surface station measurement is available and one has to rely on external data. The MeteoIO experiment employs a forcing dataset obtained through the spatial interpolation of data provided by the neighboring stations (see Sect. 3.2 and Appendix AB). GLDAS, ERA5 and ERAI experiments use different reanalysis products described in Sect. 3.3, namely GLDAS-2.1, ERA5 and ERA-Interim. Both MeteoIO and reanalysis data required to be rearranged and interpolated to 30-minute resolution in order to be used as ~~forcing~~ forcings for snow models. In the case of ERA-Interim, for example, forecasts are initialized only twice a day at 00:00 UTC and 12:00 UTC and accumulated fluxes of total precipitation, surface solar and thermal downward radiation are available as forecasts at 3-hour intervals for the following 12 hours. From these forecasts we derive the average fluxes over 3-hour intervals and we assume the fluxes to be constant during each interval. For the other ERA-Interim parameters, namely 2-meter temperature, dew-point temperature, surface pressure, 10 metre U and V wind components, we consider the analyses at 00:00, 06:00, 12:00, 18:00 UTC and the forecasts at +3 hours. These data are linearly interpolated in time to the integration time step (30 minutes) of the snow models. Some calculations are necessary to obtain all the variables required by the models. For example, ERA-Interim does not directly provide relative humidity, which we derive using the Magnus formula from the dewpoint temperature and the 2-meter air temperature (Lawrence, 2005).

The last two experiments, ERAI-LR and ERAI-BIAS, investigate ~~if it is possible to improve the performances of snow models when they are forced by reanalyses, for instance ERA-Interim, by applying two simple bias-correction methods to adjust air temperature and hence the amount of solid precipitation with respect to the total one. whether bias-correcting (some of) the reanalysis drivers improves the snow model performance. To this end, we bias-correct air temperature (and indirectly the ratio of solid to total precipitation that depends on temperature) while keeping all other variables unchanged. The idea is~~

to test whether i) the adjustment of air temperature (and the rainfall/snowfall partition) only can improve model performances and to which extent, and ii) very simple bias correction methods can be sufficient or more sophisticated ones are necessary.

In ERAI-LR we take into account the fact that ERA-Interim has a smoothed topography and the altitude of the gridpoint closest to the Torgnon station is 680 m lower than the actual elevation of the station. In ERAI-LR experiment we adjust the temperature data assuming a fixed moist lapse rate of $6.5^{\circ}\text{C}/\text{km}$. This correction consists in a cooling of 4.4°C with respect to the original temperature data. In the ERAI-BIAS experiment we correct ERA-Interim air temperature using the difference in the climatological averages between ERA-Interim data and the Torgnon station observations, which was found to be 2.1°C . This bias is assumed to be constant in time and it is subtracted from the original ERA-Interim temperature time series.

A desirable feature of each experiment is that the differences in the model outputs are mainly due to the internal model characteristics rather than to the different parameterizations used by the models to derive the solid and liquid precipitation fractions from the total precipitation input. To this end, ~~we calculate for each experiment, we estimate~~ externally the rainfall and the snowfall amounts using a fixed threshold on wet-bulb temperature. Specifically, precipitation is considered as snowfall when ~~the~~ wet-bulb temperature is lower than or equal to 1°C and as rainfall otherwise. A slightly different approach was used for GEOTOP which requires precipitation totals (rather than solid and liquid precipitation separately) and then it separates rainfall and snowfall through an internal parameterization based on a fixed threshold on dew-point temperature. In this case the dew-point temperature threshold has been calibrated to obtain approximately the same seasonal accumulated snowfall as that obtained with the method based on wet-bulb temperature. This condition is satisfied with a dew-point temperature threshold of 1.2°C . Both approaches ~~relies-rely~~ on the fact that the temperature interval where rain and snow coexists is ~~more-narrow narrower~~ for wet-bulb temperature and dew-point temperature than for air temperature, ~~thus-using~~. Using the wet-bulb or dew-point temperature allows to reduce the range for which the precipitation phase is uncertain (Sims and Liu, 2015; Endrizzi et al., 2014). With this procedure all the models are driven with the same rainfall and snowfall inputs and the differences in the model simulations are assumed to depend mainly on the model structure and on the estimated snow ablation through melting, evaporation and direct air-snow sublimation (Slater et al., 2001). This procedure is applied to the total precipitation forcing of each experiment, so also to the reanalyses, even though they provide separate snowfall and rainfall among their output variables.

5 Results

5.1 CTL - impacts of the snow model structure

We run the six models driven by the best forcing available for the Torgnon site, namely the ~~Torgnon~~-station measurements at 30-minute resolution. Figure 1 shows the simulated SWE, snow density (ρ) and snow depth (SD) time series provided by each model compared to the observations, over the period 2012-2017.

All the models are able to reproduce the overall variability of snow characteristics, although with different accuracy. The agreement between simulations and observations is evaluated in terms of centered pattern root mean square error, standard deviation and temporal correlation, and the resulting statistics are summarized through Taylor diagrams (Taylor, 2001) in Fig.

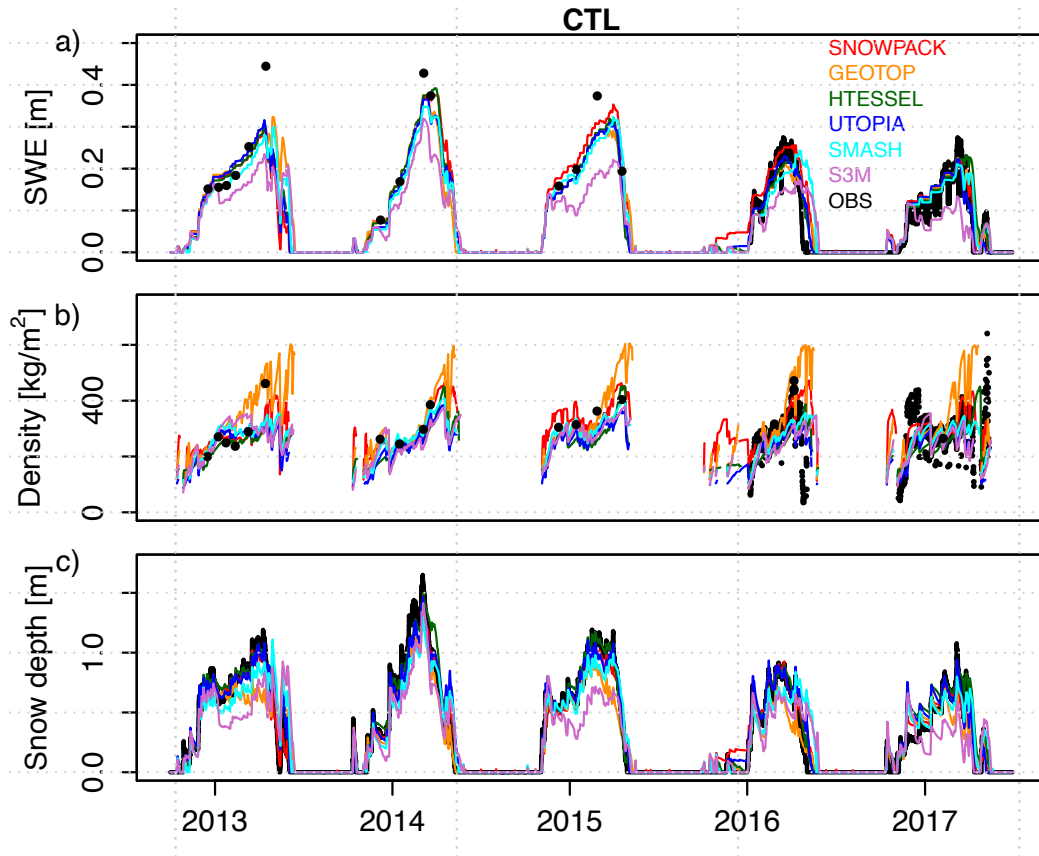


Figure 1. Results of the CTL experiment: a) snow water equivalent (SWE), b) snow density and c) snow depth simulated by the six models considered in the analysis, driven by optimal forcing, i.e. Torgnon station measurements at 30 minute resolution, over the period 2012-2017, compared to observations (black).

2. Taylor diagrams display observations as an open circle on the x-axis; the centered root mean square error between the simulated and observed variable is proportional to the distance to observations; the standard deviation of the simulated variable is proportional to the radial distance from the origin; the temporal correlation between the simulated and observed variables is shown by the angular coordinate. Evaluation metrics are calculated over simulated and observed pairs when at least one of

5 the two values exceeds a minimum threshold, namely SWE > 0.005 m, SD > 0.01 m. Snow density pairs are compared **in case both-if** the corresponding values of SWE are greater than 0.005 m. The upper panels of Fig. 2 refer to the period 2016-01-01 to 2017-06-30, when continuous measurements of all **the** three variables are available. Bottom panels refer to the full period of analysis (since 2012-10-01) for which continuous observations are available for snow depth only.

10 Snow water equivalent simulations are in **overall** good agreement with observations over the period 2016-2017 (Fig. 2a), although with some differences between the models. The best agreement is found with the SNOWPACK, HTESSEL, UTOPIA and GEOTOP models, showing the lowest errors (below 0.04 m SWE) and the highest correlations (above 0.85) with obser-

vations. SMASH and S3M are characterized by higher RMSE and lower correlation with observations with respect to the best performing models.

Snow density is simulated with lower skills compared to SWE for all models (Fig. 2b). The agreement between model simulations and observations is rather low for all models, with limited added value from highly sophisticated models. A weak correlation (lower than 0.6) and large errors (above 70 kg/m²) are found for both SNOWPACK and S3M, namely the most sophisticated and the simplest model, respectively. The GEOTOP model has clear deficiencies in representing spring snow density, in fact it exhibits an overestimation error increasing with time till the end of the snow season.

The ability of the models ~~in reproducing to reproduce~~ the temporal evolution of snow depth is related to their skills in reproducing both snow mass and density. The SNOWPACK model reproduces with high scores all the three variables, namely SWE, snow density and snow depth. In the case of GEOTOP, the overestimation of spring snow density is reflected in overall lower skills in reproducing snow depth compared to the other intermediate-complexity models (Fig. 2c). In the case of HTESSEL, instead, small errors in SWE and snow density are compensated and the model skill in reproducing snow depth ~~results slightly higher compared to is slightly higher than~~ that of the SNOWPACK model.

~~Globally the~~ The high- and intermediate-complexity models SNOWPACK, HTESSEL and UTOPIA ~~show~~ show similar and good performances in the simulation of ~~both~~ SWE and snow depth and they can be considered the best performing models. ~~Compared to them,~~ SMASH and S3M are characterized by higher RMSE and lower correlation with ~~observations, with the observations, and~~ the simplest model, S3M, ~~showing shows~~ the lowest agreement with ~~the~~ observations. In this experiment the model complexity is broadly reflected in the model performances, with the most sophisticated model performing best and the simplest model performing worst, likely owing to difficulties of the latter in representing snow melting (Fig. 1a). HTESSEL and UTOPIA, which are single-layer intermediate-complexity snow models performing almost as well as the most sophisticated model SNOWPACK, seem a good trade-off between model complexity and model accuracy when accurate meteorological forcing is employed.

We extend this analysis to a longer period of five complete snow seasons, from 2012 to 2017, limited to the snow depth variable. The relative skills of the models in reproducing snow depth over the full five-seasons period are very similar to those found for the last two-seasons period (Fig 2c,d). The RMSE of the models remains almost unchanged, while the correlation with observations slightly improves over the longer period. The behavior of the models is robust whether considering all the five seasons or only the last two seasons.

Figure 2e allows to investigate the variability of the model performances in the different snow seasons compared to the whole period. SNOWPACK, HTESSEL and UTOPIA show similar skills across different snow seasons, ~~meaning implying~~ robustness in reproducing a variety of conditions. Common simulation errors for several models are a positive SWE and a positive snow depth bias in the season 2015-2016 (Fig. 1a,c), when several ~~conditions which are challenging to simulate~~ ~~challenging conditions~~ occurred. First, ~~autumn isolated snowfalls occurred with snow-free conditions in-between~~ ~~in autumn there were isolated snowfall events separated by snowfall-free periods~~; mainly the SNOWPACK model, and to a lesser extent UTOPIA ~~model~~, failed to reproduce the rapid melting and they continued accumulating snow ~~in time~~. Second, at the end of the snow season a very rapid melting occurred, which was not captured by any of the models. All ~~the~~ models simulate a meltout date

delayed by several days with respect to the observations. Di Mauro et al. (2018) demonstrated that the accelerated snowmelt, observed in the 2015-2016 season, was caused by the deposition of mineral dust from Sahara: light absorbing impurities in snow, resulting from several dust deposition events, induce albedo reduction that alters the melting dynamics of the snowpack hence advancing favouring snowmelt. As none of the model-models used in this study accounts for the impact of impurities on snow dynamics ~~-,modeled-~~(and in any case no information on dust deposition is provided to the models) simulated snow melt dates in 2016 were~~not-surprisingly-, not surprisingly,~~ significantly delayed.

The GEOTOP, SMASH and S3M models show different skills depending on the snow season (Fig. 2e) and they provide a wider range of variability in their agreement with the observations compared to SNOWPACK, HTESSEL and UTOPIA. For example, a season which is relatively simple to reproduce by all models is 2013-2014. An abundant but ephemeral snow cover was properly accumulated and melted by all models. After a snow-free period, the onset of a persistent snow cover was sustained by heavy snowfalls which led to the highest snow peak in the study period. After this peak, the melting has been quite steady, with few spring snowfall events. These conditions allow all models, even the simplest one, to accurately reproduce the snowpack evolution in terms of snow mass and depth. As a result, for this season the differences between the models in terms of RMSE, standard deviation and temporal correlation with observations are smaller than for other seasons. On the contrary, the season 2012-2013 is more difficult to reproduce for some models, namely GEOTOP, SMASH and S3M, than for SNOWPACK, HTESSEL and UTOPIA. This season was characterized by many snowfall episodes of moderate and light intensity, with moderate melting in-between. In the second half of May 2013 a series of late snowfalls restored a temporary snow cover with more than 0.5 m depth that gradually melted in a couple of weeks. In this conditions, SNOWPACK, HTESSEL and UTOPIA are able to ~~represent quite accurately~~accurately represent the changes in the snow depth, while GEOTOP, SMASH and S3M generally tend to overestimate snow depth.

GEOTOP ~~sistematically~~systematically overestimates snow density with increasing errors from late winter to the end of the snow season. These errors are reflected in the snow depth simulations: spring snow depth and the snow depth peak are underestimated in each snow season of the study period. SMASH, for the 2012-2013 ~~season-as-well-as-for-and~~ the 2015-2016 ~~season, misrepresents the correct seasons, delays the~~ timing of the snow depth and snow mass peaks~~and delays them in the snow season~~. The delay in the representation of the snow peaks is almost fully compensated by ~~a too~~an excessively rapid spring melting which allows to keep the date of ablation relatively close to the observed one. S3M ~~sistematically~~systematically underestimates both snow depth and snow water equivalent during all the snow seasons, while the snow density is within the range of variability of the model ensemble. It follows that for S3M the critical variable to improve is SWE.

In conclusion, an added value of sophisticated and intermediate-complexity models compared to lower-complexity models emerges especially during snow seasons that have a more complex temporal behavior.

5.2 RAD-ERA1, SWIN-CLS - model sensitivity to the radiation input

A typical problem occurring in case of snowfall is that when the radiation sensors get covered ~~by~~with snow they record inaccurate ~~or even wrong~~ data. To take into account this issue and test how it affects snow simulations, in the experiment RAD-ERA1 we use incoming longwave and shortwave radiation data from the Torgnon station except in case of snowfall, when we

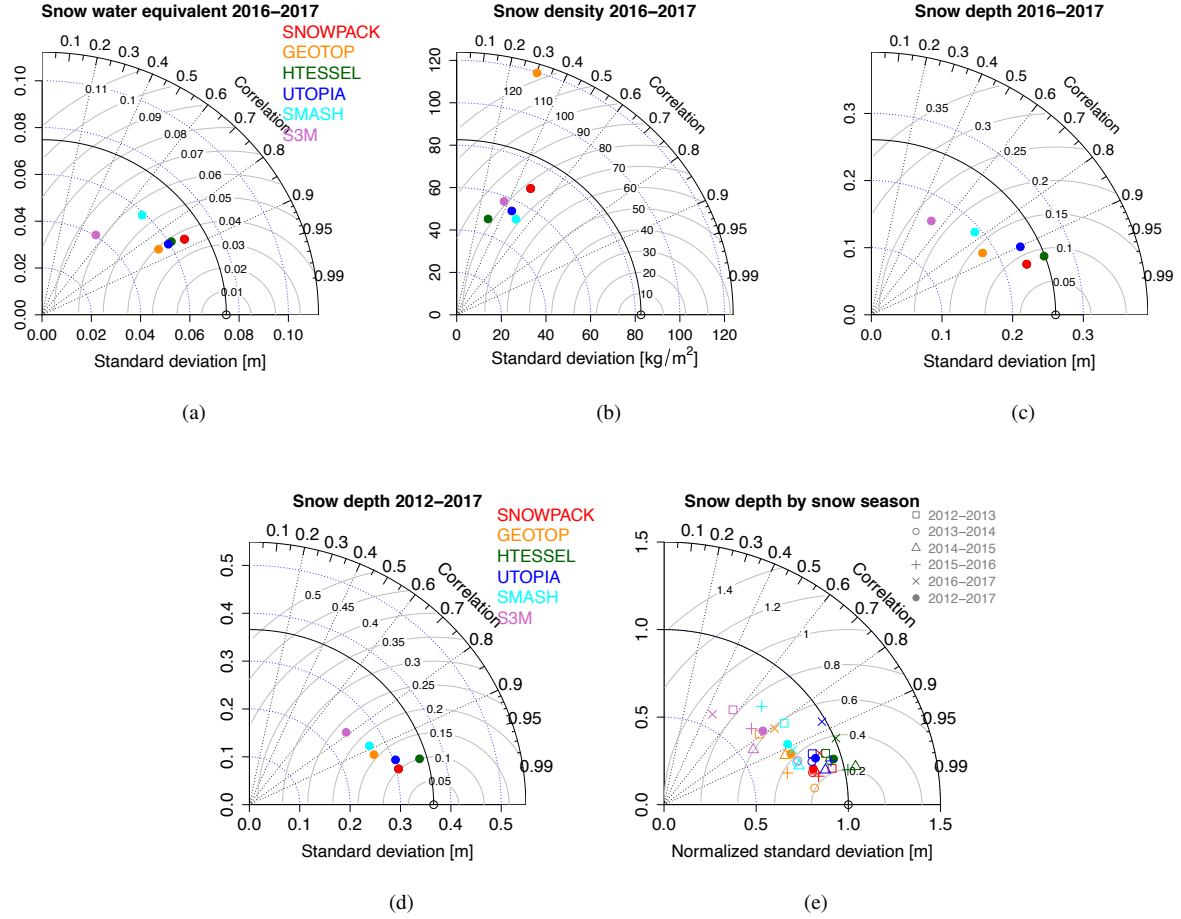


Figure 2. Taylor diagrams of the modeled vs. observed a) snow water equivalent, b) snow density and c) snow depth in the control experiment (CTL) for the period 2016-01-01 to 2017-06-30. Bottom panels represent the statistics of snow depth d) for the whole period 2012-2017 and e) for each snow season in the same period. **Please note that, differently** from other panels, in panel e) model standard deviations are normalized with respect to the observed ones.

Table 3. RMSE, bias and Pearson correlation of snow depth simulations with respect to observations for each model and each snow season in the control experiment (CTL).

RMSE							
Model	2012-2013	2013-2014	2014-2015	2015-2016	2016-2017	Avg.	Std.
SNOWPACK	0.08	0.12	0.10	0.12	0.08	0.10	0.02
GEOTOP	0.21	0.14	0.18	0.13	0.15	0.16	0.03
HTESSEL	0.11	0.13	0.07	0.08	0.15	0.11	0.03
UTOPIA	0.11	0.14	0.07	0.12	0.13	0.11	0.03
SMASH	0.19	0.19	0.12	0.23	0.12	0.17	0.05
S3M	0.28	0.20	0.28	0.21	0.24	0.24	0.04
BIAS							
Model	2012-2013	2013-2014	2014-2015	2015-2016	2016-2017	Avg.	Std.
SNOWPACK	-0.04	-0.03	-0.05	0.10	0.01	0.00	0.06
GEOTOP	-0.05	-0.11	-0.11	-0.06	-0.01	-0.07	0.04
HTESSEL	0.04	0.04	0.01	0.06	0.11	0.05	0.04
UTOPIA	0.03	-0.01	-0.01	0.09	0.06	0.03	0.04
SMASH	-0.04	-0.09	-0.04	0.05	0.04	-0.02	0.06
S3M	-0.08	-0.12	-0.21	0.01	-0.10	-0.10	0.08
Pearson Correlation							
Model	2012-2013	2013-2014	2014-2015	2015-2016	2016-2017	Avg.	Std.
SNOWPACK	0.98	0.98	0.97	0.98	0.94	0.97	0.01
GEOTOP	0.79	0.99	0.92	0.97	0.81	0.90	0.09
HTESSEL	0.95	0.96	0.98	0.98	0.93	0.96	0.02
UTOPIA	0.94	0.96	0.98	0.96	0.87	0.94	0.04
SMASH	0.81	0.95	0.96	0.69	0.91	0.86	0.11
S3M	0.57	0.95	0.84	0.74	0.45	0.71	0.20

employ external LWIN and SWIN data derived from ERA-Interim. In the other experiment, SWIN-CLS, we replace observed incoming shortwave radiation data with the external data described in Sect. 4. The results of these simulations are reported in Table 4. Although the difference between external data and Torgnon data can be high at the time step of the model (not shown), their overall impact on snow simulations is low. In fact, for each model we obtain values of RMSE close to those obtained
5 in the CTL experiment. In particular, model skills ~~does~~do not improve using ERA-Interim or interpolated incoming radiation forcing.

5.3 TIME-3h, TIME-6h, TIME-12h - model sensitivity to the temporal resolution of the forcing

A common condition when modeling snowpack evolution in data-sparse areas is the unavailability of meteorological forcings with high temporal resolution, as high as 30 minutes, like that employed in the CTL experiment. In this section we assess the sensitivity of the models to the temporal resolution of the forcing. To this aim, the original ~~30-minute-resolution~~ meteorological observations at Torgnon, with 30-minute resolution, have been sampled every 3, 6 and 12 hours, and then linearly interpolated at ~~a the~~ finer (30-minute) time step, with the only exception of total precipitation that has been accumulated over the 3, 6, 12 hour periods and then equally distributed among the 30-minute sub-periods. Incoming shortwave radiation for the TIME-12h experiment has been derived with two different methods, i.e. by linear interpolating the measurements at 00:00 and 12:00 and by rescaling the potential radiation at 30 minute temporal resolution to the observed radiation at 12:00 (see Sect. 4 for details).

As expected, the longer the sampling period the smoother are the input time series. For these three (and the other remaining six) experiments, we show in Fig. 3 the biases of air temperature, total precipitation, rainfall and snowfall forcings with respect to the reference forcing of the CTL experiment. Given the method employed to derive TIME-3h, TIME-6h, TIME-12h forcings we expect no bias for total precipitation, while some differences can arise in the rainfall/snowfall partition owing to possible differences in air temperature. According to Fig. 3, TIME-3h and TIME-6h air temperature biases are close to zero, while TIME-12h air temperature bias is about 0.5°C , with the effect of reducing the amount of the solid precipitation by 10%. We investigate the impact of these biases on the snow simulations in the following.

Figure 4 represents, for all the models, the simulated snow depth and SWE when input data are sampled (or accumulated, in case of total precipitation) at 3, 6 and 12 hours and then interpolated (or equally distributed for precipitation) over 30-minute time steps, compared to the simulated snow depth obtained with the original 30 minute resolution forcing (CTL) and compared to observations. The TIME-12h experiment employs the incoming shortwave radiation estimated with the potential radiation method; Appendix C reports the corresponding experiment employing the incoming shortwave radiation estimated with linear interpolation of the station measurements. In addition, Table 4 reports the RMSE associated with all these simulations.

The model response to the degradation of the temporal resolution of the forcing ~~can vary remarkably depending depends~~ on the model and season. A common feature of the models is the small (or null) difference in terms of RMSE between TIME-3h and CTL simulations, indicating that using forcing data at 3-hour temporal resolution generates snow depth simulations almost as accurate as in the case of 30-minute resolution input data.

A second common feature is the general worsening of model performances when using input data at ~~a lower temporal resolution than 3 hours~~ 6-hour temporal resolution, reflected into an increase of the RMSE values. TIME-6h simulations are usually very close to the CTL in the accumulation period up to the snow peak. Afterwards, in the melting period, some models, mainly SNOWPACK, UTOPIA and to a lesser extent HTESSEL and SMASH, slightly overestimate snow mass/depth in selected seasons, contributing to an increase in the model RMSE. ~~The largest RMSE values are generally obtained with input data at 12 hour resolution: in this case all models represent a thinner snow depth throughout the snow season compared to the CTL run. The snow depth underestimation in~~ Compared to the TIME-6h experiment, the TIME-12h ~~experiement can be related to biases in the forcing, namely the positive air temperature bias and the underestimation of snowfall by 10%~~

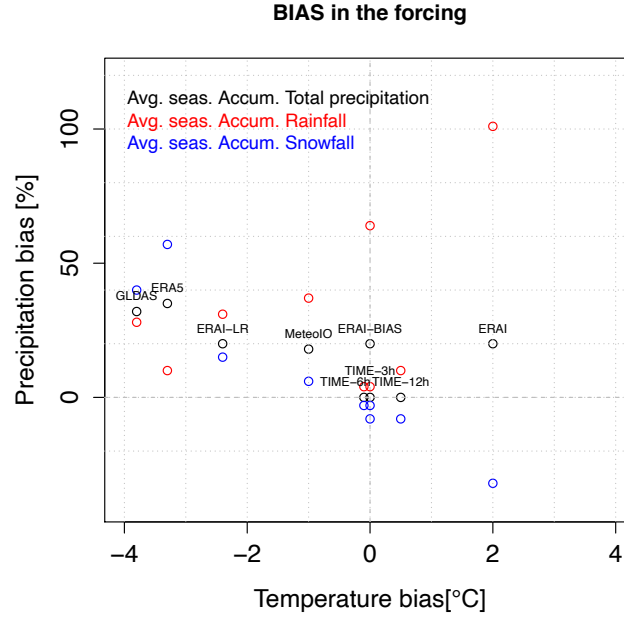


Figure 3. Temperature, total precipitation, rainfall and snowfall average seasonal biases in the forcings employed in each experiment with respect to the Torgnon station measurements.

(Fig. 3) combined with a poor representation of the incoming shortwave radiation. In fact, TIME-12h incoming shortwave radiation forcing exceeds, in average, the CTL values by $+97 \text{ W/m}^2$. This large experiment with incoming shortwave radiation interpolated with the linear method (TIME-12h-LIN) shows higher RMSE on snow depth and a clear worsening of model performances (Table 4). In the TIME-12h-LIN experiment the overestimation of the incoming shortwave radiation explains most of the (see Section 4) causes an underestimation of the surface snow depth. On the contrary, the TIME-6h incoming shortwave radiation forcing shows, a small negative bias (-7 W/m^2) compared to the CTL forcing, which can contribute to the overestimation of the snow depth at the end of the snow season. TIME-12h experiment with SWIN estimated with the potential radiation method (TIME-12h-SWINPOT) shows improved model performances compared to both the TIME-12h-LIN and TIME-6h experiments for SNOWPACK, HTESSSEL and UTOPIA, with the former two models showing RMSE comparable with that of the CTL run. GEOTOP and S3M show similar skills in the TIME-12h experiments and higher RMSE in the TIME-12h experiments compared to the TIME-6h experiment and the CTL run. Finally the SMASH model shows little or no differences between the TIME-12h experiments and the TIME-6h, TIME-3h and CTL experiments.

The different sensitivities of the models In conclusion, the six models show different sensitivities to the bias in the forcing is quite impressive. While for models with higher complexity, Models with high- and intermediate-complexity (SNOWPACK, GEOTOP, HTESSSEL and UTOPIA) the error on snow depth increases at increasing temporal are sensitive to both the time degradation of the forcing, for simpler models (SMASH and and to the method used to interpolate the 12-hourly SWIN.

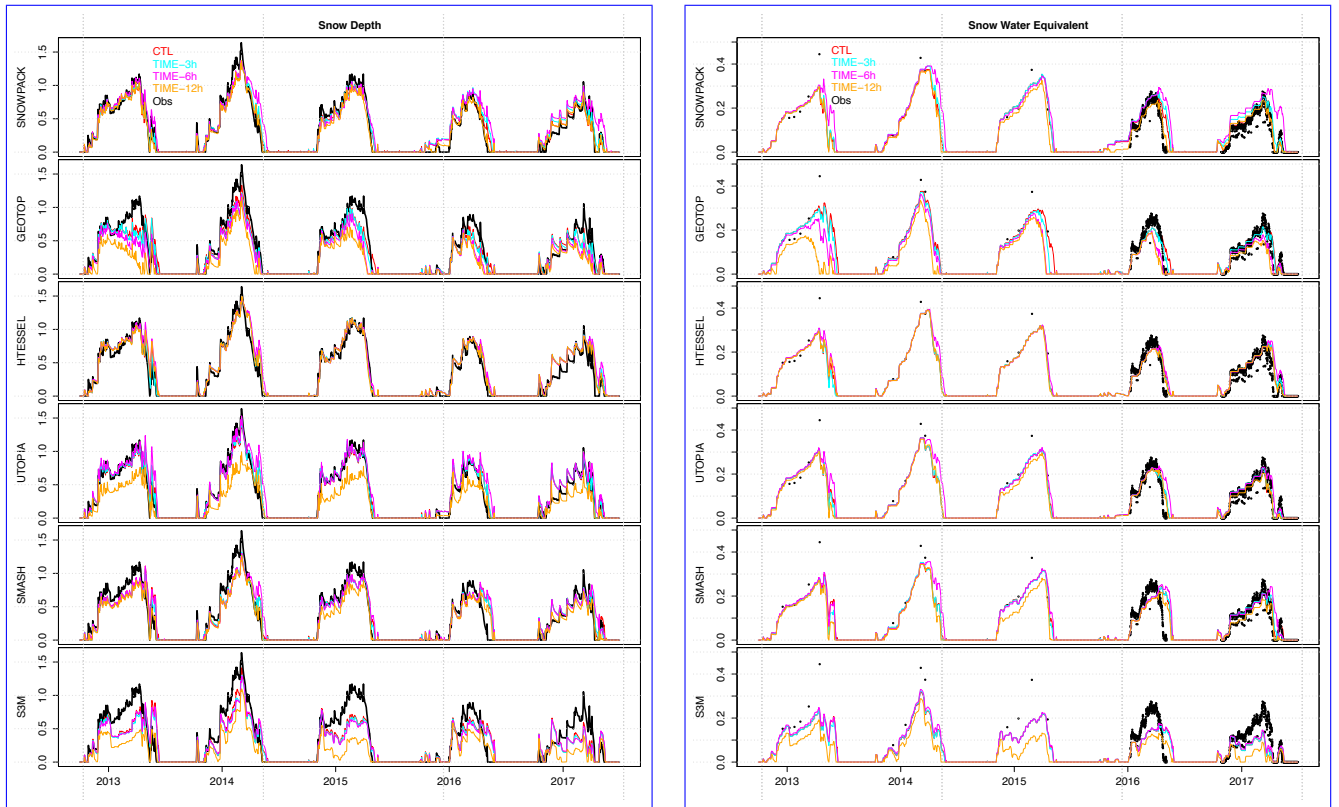


Figure 4. Model simulations of snow depth and SWE when the input is sampled at 3, 6 and 12 hours, compared to the CTL run and to the observations. The TIME-12h experiment employs the incoming shortwave radiation estimated with the potential radiation method.

GEOTOP and S3M) the performance is less affected. SMASH, for example, shows little differences between TIME-3h, TIME-6h, TIME-12h and the CTL simulations, revealing to be the less sensitive model are sensitive to the time resolution degradation of the forcing . In the TIME-12h experiment, SMASH provides the smallest error on snow depth and a similar error on SWE as the most sophisticated model, SNOWPACK. Similarly to SMASH , S3M shows very little differences between TIME-3h, TIME-6h simulations and the CTL simulation, while its performances in the TIME-12h experiment falls within the range of values obtained for intermediate-complexity models. In our experiment the lower-complexity models reveal a lower sensitivity to errors in the meteorological forcing , and when these models are driven with 6-hourly but not to the method used to interpolate the 12-hourly SWIN, and finally SMASH shows low sensitivity to both the time degradation of the forcing and to the method used to interpolate the 12-hourly input data they provide similar performances as selected intermediate-complexity models. So, with low temporal-resolution data the performances of intermediate and low complexity models get comparable. However, from SWIN.

From these experiments an added value of the most sophisticated model SNOWPACK emerges. SNOWPACK forced by the 12-hour resolution forcing still provides lower errors than the simplest model S3M forced by the best available forcing at 30 minute temporal resolution (Table 4).

Concerning the TIME-12h experiment, the SWIN forcing derived with the potential radiation method provides overall better results compared that derived from linear interpolation of the station measurements. In the following we will further analyze the TIME-12h experiment with SWIN forcing derived with the potential radiation method.

5.4 MeteoIO, GLDAS, ERA5 and ERA-Interim - model sensitivity to the spatial resolution and bias in the forcing

We consider a rather standard case for which no station measurements are available for the area of interest and one has to rely on gridded datasets, which are generally characterized by lower resolution and lower accuracy with respect to station measurements. To explore a representative range of possible alternatives we employ datasets with different characteristics: the MeteoIO dataset, based on the interpolation of data from neighboring stations, GLDAS, ERA5 and ERA-Interim reanalyses at 25, 31 and 80 km respectively. An overview of the comparison between the meteorological forcing provided by these datasets and the observations in Torgnon is shown in Fig. 3.

The MeteoIO forcing is in fairly good agreement with observations. Compared to the meteorological measurements at the Torgnon station, MeteoIO shows an average bias of -1°C per snow season and about 20% overestimation of the seasonal total precipitation. However, the effect of these biases on the solid precipitation is weak, and the average seasonal snowfall is very close to the observations. When the MeteoIO forcing is used, the best agreement between simulated and observed SWE and snow depth is obtained with the GEOTOP and SMASH models. Both models provide similar RMSE values as in the CTL runs. The S3M model exhibits a moderate decrease in the model performance when driven by MeteoIO compared to CTL, with lower RMSE than the HTESSEL and UTOPIA models. Conversely, the SNOWPACK, HTESSEL and UTOPIA model errors are respectively more than twice and three times the corresponding errors in the CTL run, revealing a remarkable sensitivity to relatively small errors in the input data. The S3M model exhibits a moderate decrease in. Despite a relatively small average error in the temperature input (-1°C), the daily differences are generally stronger in winter and they can reach values exceeding -4°C . The main issue in snow model simulations is the overestimation of snow depth in winter (in selected snow seasons) and in spring (always). A plausible explanation for these errors is that colder-than-observed winter temperatures might favor the development of a cold snowpack which melts too slowly. Consequently, the model performance when driven by MeteoIO, with lower RMSE than the HTESSEL and UTOPIA models models tend to overestimate the snow at the surface and to predict a delayed ablation date.

The GLDAS forcing is affected by strong cold and wet biases. With an average temperature bias of -3.8°C and an average a strong cold bias, with average temperature differences of -3.8°C compared to the observations, and by a moderate total precipitation bias of 296%, GLDAS lies outside the range displayed in +32% in average over the considered seasons (Fig. 3). As expected, the large errors in GLDAS forcings lead to huge errors in the GLDAS temperature forcing lead to large errors in the simulated snow water equivalent and depth for all models, as confirmed by RMSEs lying out of range in in Table 4 and Fig. 6a,c,d. The magnitude of the error on snow depth shows large variations from season to season: snow depth estimates

Table 4. RMSE values associated to snow depth and snow water equivalent simulations for all models and all experiments over the periods 2012-2017 and 2016-2017, respectively.

Exp	RMSE snow depth [m]					
	SNOWPACK	GEOTOP	HTESSEL	UTOPIA	SMASH	S3M
CTL	0.10	0.17	0.11	0.12	0.17	0.25
RAD-ERA1	0.12	0.17	0.14	0.13	0.17	0.25
SWIN-CLS	0.11	0.21	0.12	0.13	0.18	0.24
TIME-3h	0.12	0.19	0.11	0.12	0.16	0.26
TIME-6h	0.17	0.26	0.15	0.18	0.19	0.27
TIME-12h(SWINPOT)	0.21(0.11)	0.37(0.35)	0.44(0.12)	0.38(0.26)	0.17(0.17)	0.38(0.39)
MeteoIO	0.23	0.20	0.38	0.40	0.19	0.31
GLDAS	1.99 0.67	2.45 0.41	3.61 0.79	3.46 0.49	1.89 0.63	3.55 0.84
ERA5	0.74	0.34	0.76	0.80	0.71	0.85
ERA1	0.18	0.45	0.20	0.20	0.27	0.32
ERA1-LR	0.54	0.20	0.58	0.67	0.36	0.46
ERA1-BIAS	0.18	0.27	0.20	0.26	0.13	0.16

Exp	RMSE SWE [m]					
	SNOWPACK	GEOTOP	HTESSEL	UTOPIA	SMASH	S3M
CTL	0.04	0.04	0.04	0.04	0.06	0.08
RAD-ERA1	0.06	0.04	0.05	0.04	0.06	0.08
SWIN-CLS	0.05	0.04	0.04	0.03	0.06	0.07
TIME-3h	0.06	0.03	0.04	0.03	0.06	0.08
TIME-6h	0.09	0.05	0.05	0.05	0.07	0.07
TIME-12h(SWINPOT)	0.05(0.03)	0.07(0.07)	0.13(0.04)	0.13(0.03)	0.05(0.05)	0.10(0.10)
MeteoIO	0.10	0.04	0.13	0.13	0.07	0.11
GLDAS	1.31 0.38	1.68 0.22	1.77 0.38	1.82 0.16	1.19 0.33	1.56 0.38
ERA5	0.28	0.12	0.22	0.24	0.26	0.24
ERA1	0.05	0.12	0.05	0.05	0.08	0.08
ERA1-LR	0.19	0.04	0.18	0.19	0.13	0.15
ERA1-BIAS	0.05	0.05	0.03	0.05	0.03	0.05

are relatively close to the observations in the first three snow seasons while a large overestimation occurs in the last two snow seasons. This behavior can be linked to the error on the total precipitation, up to +129% and +102% relative to observations in the last two snow seasons.

ERA5 has a large temperature bias (-3.3°C) ~~like GLDAS, but a lower~~ and a moderate precipitation bias (+35%), similarly to GLDAS. The combined effect on the snowfall input is an excess of more than 50% compared to observations (Fig. 3), which clearly affects the snow model output. As expected, all models overestimate snow depth and the duration of the snow cover. The models tend to reproduce a similar evolution of snow depth as in the CTL experiment but with thicker snowpack. In detail, SNOWPACK, HTESSEL and UTOPIA give similar snow depth outputs, consistently with the behavior found in the CTL run. GEOTOP provides the lowest RMSEs for snow water equivalent and snow depth, but this is mainly due to a compensation between the error in the ERA5 forcing (leading to overestimation) and the model error identified in the CTL experiment (leading to underestimation). In general, the difference in performance between models of different complexity is reduced when the ERA5 forcing is used. For example the RMSE is similar for the simplest model SMASH and the most sophisticated model SNOWPACK, as it is for S3M and HTESSEL or UTOPIA.

The ERA-Interim forcing (ERA-Interim) shows a $+2^{\circ}\text{C}$ temperature bias and a snowfall deficit of about 30% compared to the Torgnon observations. When forced by ERA-Interim data, GEOTOP, SMASH and S3M underestimate snow depth in all seasons, while SNOWPACK, HTESSEL and UTOPIA underestimate snow depth mainly during the season 2014-15, when the ERA-Interim snowfall is considerably lower ~~with respect to~~ than the observations throughout this snow season (Fig. 5a). In other snow seasons, for example 2013-14, 2015-16 and 2016-17, SNOWPACK, HTESSEL and UTOPIA snow depth simulations are in fairly good agreement with the observations (see for example Fig. 5b). Overall, SNOWPACK, HTESSEL and UTOPIA provide relatively good results when forced by ERA-Interim, with a moderate loss of accuracy with respect to the case of optimal forcing (CTL). In the following we explore the possibility to reduce the RMSE of the other intermediate- and low-complexity models by correcting the main biases in the meteorological forcings.

5.5 Impact of the bias adjustment of ERA-Interim air temperatures

We test the effect of two very simple bias correction techniques applied to the ERA-Interim air temperature. In the first approach, in the ERA-LR experiment, we take into account the difference in elevation between the ERA-Interim gridpoint at Torgnon and the true elevation of this site, applying a lapse rate correction, i.e. subtracting 4.4°C from the original ERA-Interim data. Alternatively, in the ERA-Interim-BIAS experiment, we remove the average bias of ERA-Interim data at the Torgnon site with respect to the station measurements, i.e. subtracting 2.6°C from the original ERA-Interim data.

The lapse rate correction ~~reduces too much~~ excessively reduces ERA-Interim temperatures, ~~in fact~~: the average temperature bias shifts from +2 to -2.4°C and the snowfall amount increases from -32 to +15% (Fig. 3).

The net effect on the model outputs (ERA-Interim-BIAS experiment) is an overestimation of snow water equivalent and snow depth. With respect to the ERA-Interim experiment, the RMSE values increase for all models except for GEOTOP, which actually shows a good agreement with observations during the seasons 2013-14 and 2015-16, while it overestimates snow depth in the first

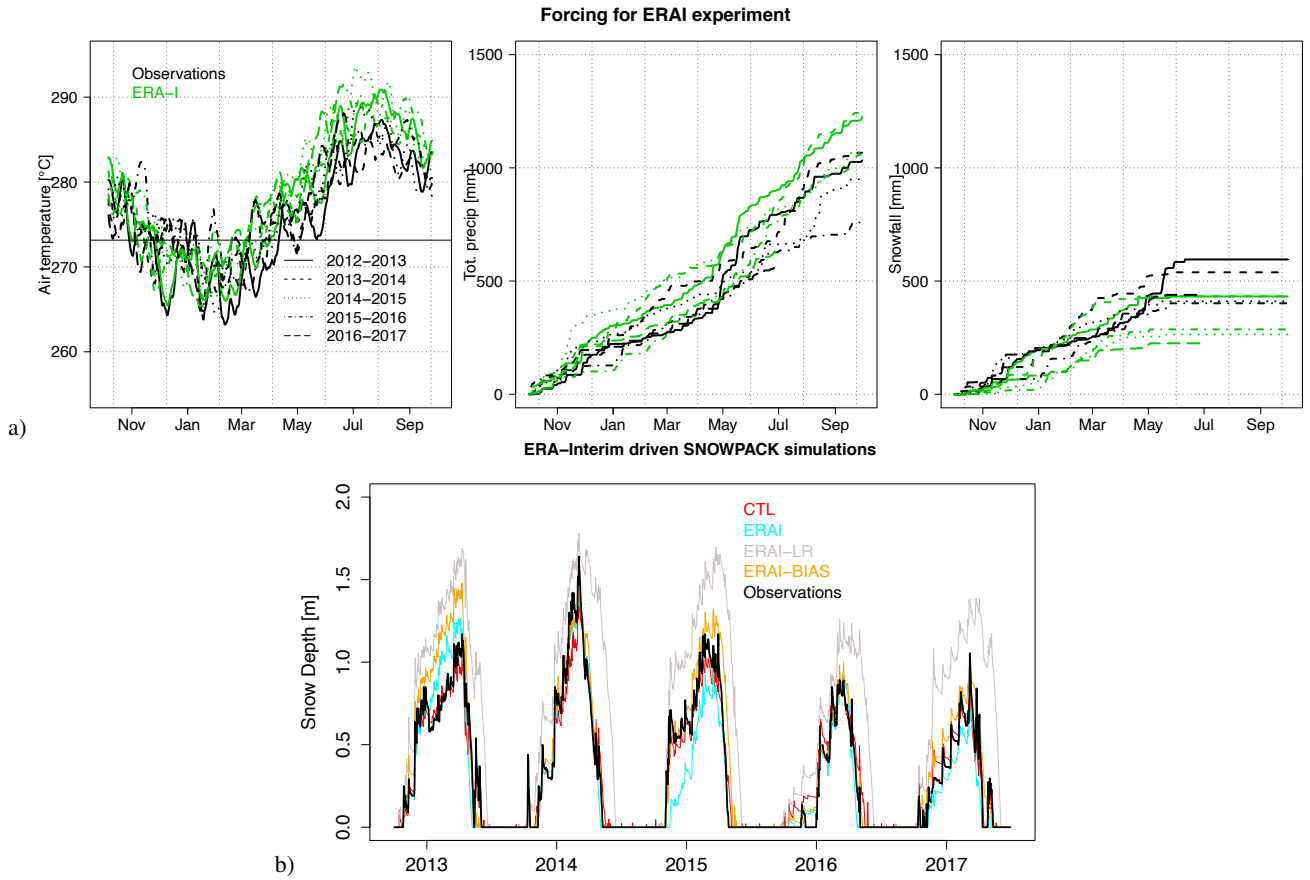


Figure 5. a) ERA-Interim air temperature, total precipitation and snowfall (derived as explained in Sect. 4) at the Torgnon site, compared to the station measurements (black) for each snow season of the period 2012-2017; b) Example of ERA-Interim driven snow depth simulations (ERAI, ERA-LR and ERAI-BIAS experiments), obtained ~~with using~~ the SNOWPACK model, compared to CTL run and snow depth observations.

half of the other seasons. The GEOTOP underestimation error observed in the ERAI experiment is ~~here-compensated-by-too~~ compensated by excessively cold input air temperature, which ~~favor~~ favor-favors the development and ~~the~~ duration of the snowpack.

The correction based on the adjustment of the mean ERAI temperature bias (ERAI-BIAS experiment) allows to almost remove the snowfall bias. Therefore, this approach guarantees the most effective correction to improve the agreement of the forcing data with the Torgnon station measurements. Clearly this approach requires to know at least the average temperature at the site of interest. This correction successfully reduces the RMSE on snow water equivalent and snow depth simulations with respect to the corresponding runs driven by the raw ERA-Interim data for GEOTOP, SMASH and S3M ~~models~~. For the most sophisticated SNOWPACK model, the correction applied to ERA-Interim data has no effects on the RMSE values of snow water equivalent and snow depth simulations, which remain unchanged. ~~In fact, while~~ While the simulated snow depth is generally close to observations, improvements gained in the selected seasons (i.e. 2014-15) are compensated by lower performances in

other (2012-13) seasons (Fig. 5d), so that, in average, the overall effect on the RMSE is negligible. For the UTOPIA model, the correction applied to ERA-Interim data has no effects on the snow water equivalent, however it slightly increases the error on snow density (lower correlation with available observations) and thus the error on snow depth simulations.

5.6 Discussion

5 While much work has been done to characterize the performances of snow models when driven by accurate input data (e.g. Vionnet et al., 2012; Boone and Etchevers, 2001; Bartelt and Lehning, 2002; Dutra et al., 2010), model responses depend-
 ing on different degrees of accuracy of the input data still needs to be explored in detail. This study ~~offers some hints~~ sheds
light on this research topic by assessing the simulations of six state-of-art snow models driven by input data with varying
 accuracy, focusing on the fully-instrumented Torgnon site, in the NW Italian Alps. The snow models selected for the analysis
 10 ~~present also are characterized by~~ different degrees of complexity, from highly sophisticated multi-layer snow models to rather
 simple single-layer ~~snow~~ models, with the aim of ~~shedding light also on the~~ exploring relations and trade-offs between model
 complexity and model performances in reproducing snowpack dynamics.

In our experiment, in the case of optimal forcing, namely Torgnon station data at 30-minute resolution, the most sophisticated
 model SNOWPACK and the intermediate-complexity models HTESSEL and UTOPIA show the best agreement with observa-
 15 tions. In particular HTESSEL and UTOPIA ~~models~~, with their single-layer, simpler snow schemes compared to SNOWPACK,
 can be considered a good trade-off between model complexity and model accuracy. When considering snow depth simulations,
 for which validation data are available for a longer period than for SWE, an added value of these high- and intermediate-
 complexity models compared to lower complexity models is evident, especially in the snow seasons that are more difficult
 to reproduce. SNOWPACK, HTESSEL and UTOPIA ~~, in fact,~~ show similar and good performances across different seasons,
 20 revealing robustness in reproducing a variety of conditions, while the simpler snow models SMASH and S3M show larger
 dispersion of the seasonal scores.

Snow density is more difficult to simulate than SWE and snow depth for all models. The correlation between model sim-
 ulations and observations is quite low for all models, with no clear added value from highly sophisticated ones (Fig. 2b). ~~In~~
~~fact, similar RMSE values with respect to snow density observations are found for SNOWPACK, UTOPIA, SMASH and S3M~~
 25 ~~models; GEOTOP, instead,~~ GEOTOP provides a much larger error compared to the other models, especially in the spring
 season, suggesting further checks on the snow density parameterization.

The response of the snow models forced by gradually lower accuracy data is summarized in Fig. 6, showing the model
 RMSE for all experiments and all variables (upper panels) and the complementary information on the model ranking (bottom
 panels). No remarkable differences can be detected in the model skills when using alternative radiation data instead of the
 30 Torgnon station measurements, as done in experiments RAD-ERA-Interim and SWIN-CLS. The substantially equivalent results ob-
 tained replacing measured data with ERA-Interim data in case of snowfall (RAD-ERA-Interim experiment) can be explained by the
 combination of two conditions: the intermediate elevation of 2160 m a.s.l. and the orientation of the Torgnon site, both likely
 contributing to a rapid melting of the snow obstructing the radiometer. This adjustment does not affect model performances.
 Similar results are found employing ~~SW~~ SWIN radiation estimated as clear-sky radiation attenuated by a factor based on MSG

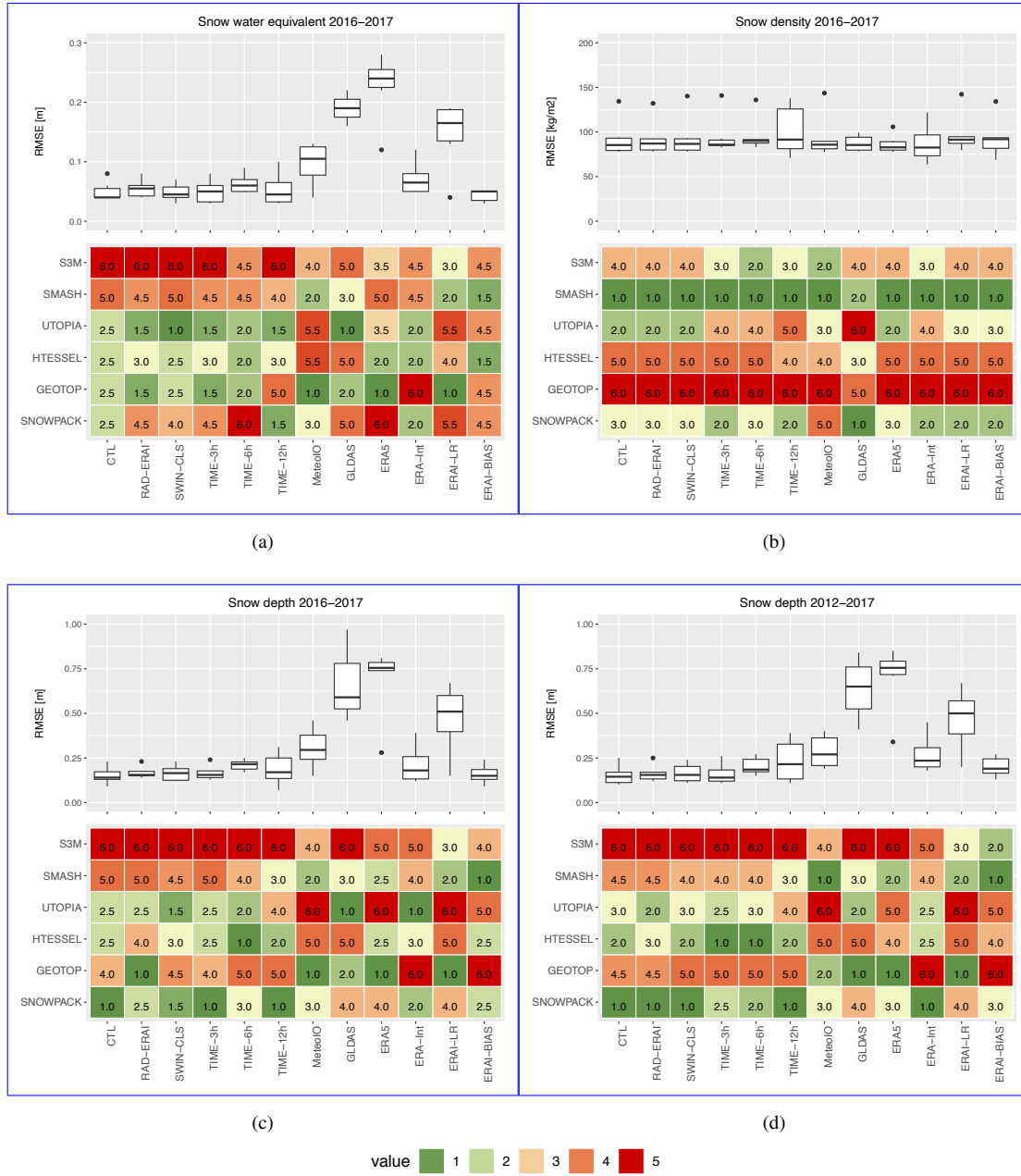


Figure 6. Root Mean Square Error associated to a) snow water equivalent, b) snow density and c) snow depth simulations for each experiment and each model over the period 2016-2017. Panel d) shows the same statistics as c) but on the whole period 2012-2017. **Values for GLDAS experiments lie outside the plotted range.** Upper panels represent the boxplot statistics, lower panels represent the model rank (1=best model, 5=worst model).

cloud mask and neighbouring station radiation measurements (SWIN-CLS, Sect. 4). Each model shows ~~almost identical RMSE~~ in similar RMSE on snow depth in the CTL, RAD-ERA-Interim and SWIN-CLS experiments.

The use of accurate meteorological inputs but at lower temporal resolution, for instance Torgnon station data sampled at 3 hourly time step and then interpolated to the model time step, does not affect model performances ~~in our experiment~~. At longer time steps we observe a gradual decay of the snow model skills, which is strongest at the longest time step considered, corresponding to 12 hours. Therefore Similar results were obtained in a previous study in which the original forcing was averaged in time over 3 hours and the resulting time series was interpolated to the model time step (Ménard et al., 2015). Therefore we can conclude that, at their typical temporal resolution (3 hours), climate and weather forecast model simulations, as well as reanalysis data, can be suitable for driving snowpack models, ~~while lower temporal resolution input data require~~. The use of input data with temporal resolution lower than 3 hours requires more in-depth consideration. ~~For as we observe a gradual decay of the snow model skills for most models. With 12-hourly resolution input, for example, the sampling of the original 30-minute forcing at 12 hours implied a worsening of the shortwave incoming radiation input, that could be partly overcome with more sophisticated interpolation techniques than the basic linear interpolation employed in this study. In our experiment, when incoming shortwave radiation is found to be a key variable affecting the model performances. While the simple linear interpolation of the 12-hourly forcing is employed, the intermediate-complexity model SMASH provide comparable or lower errors than any other more sophisticated model. This suggests that with low temporal resolution forcing selected intermediate-complexity models can be employed without reducing the accuracy of the snow outputs.~~

When only low temporal resolution forcing is available, more sophisticated interpolation techniques compared to the basic linear interpolation employed in this study could improve the agreement with the reference data. When employing low temporal resolution forcing, similar performances are found for intermediate- and low-complexity models, probably because for the low-complexity models, the model error and radiation data to the model time step provides poor SWIN estimates and poor snow model performances, a slightly more sophisticated method based on the scaling of the potential radiation on the SWIN measurements at 12:00 allows to improve the snow simulations and to obtain model skills comparable to or even better than the TIME-6h experiment. With this second method the bias on the incoming shortwave radiation flux is almost completely cancelled out. A residual negative bias (-7 W/m^2) of the incoming shortwave radiation in the forcing error are partly compensated. TIME-6h experiment contributes to the overestimation of the snow depth at the end of the snow season. For SNOWPACK and HTESSEL the 12-hourly forcing with improved SWIN input allows to obtain surprisingly good performances, as shown by the comparable RMSEs in the TIME-12h experiment and in the CTL run.

Where meteorological station data are not available, spatial interpolation of neighboring stations data or reanalyses can be a valid alternative. In our experiment the best results are obtained with ERA-Interim forcing. ~~Indeed, despite~~ Despite the coarse spatial resolution, ERA-Interim satisfactorily reproduces the meteorological conditions at the Torgnon gridpoint (Fig. 3) and the model errors in terms of snow depth and snow water equivalent are only slightly higher than in the CTL experiment (Fig. 6a,c,d). SNOWPACK, HTESSEL and UTOPIA ~~models~~ again provide the lowest errors compared to intermediate- and low-complexity snow models (GEOTOP, SMASH, S3M). However, also the latter can be an interesting option after applying a simple adjustment of the average ERA-Interim temperature bias with respect to the Torgnon station data, and consequently

adjusting also the snowfall amount. In ~~fact, in~~ this way the performances of the intermediate- and lower complexity snow models (GEOTOP, SMASH, S3M) can be substantially improved. The temperature adjustment based on the lapse rate (ERAI-LR), accounting for the difference in elevation between the ERA-Interim gridpoint and the real elevation of the Torgnon station, is found to worsen the model performance. In fact, this correction is blind to the local climatic features and might not be suitable in all situations. For example, in this case the lapse-rate correction is too large and it causes a temperature bias of similar amplitude but opposite sign with respect to the original ERA-Interim data. As a general remark, it is preferable to apply a temperature correction based on local temperature observations or even just climatology, when available, as the correction based on the lapse rate does not ensure a better agreement with the reference data.

Spatial interpolations of neighboring station data, such as the MeteoIO interpolation used here, can be another valid alternative in ~~case of~~ absence of in-situ observations. In our experiment the models RMSE values for snow water equivalent and snow depth are generally ~~comparable to those obtained with~~ higher than those obtained using the Torgnon data at ~~12-hourly~~ lower temporal resolution. GLDAS ~~, and to a lesser extent ERA5, are affected by large temperature and precipitation biases~~ ERA5 are affected, on average, by a large temperature bias and a moderate precipitation bias at the Torgnon gridpoint, probably owing to difficulties of these datasets in simulating processes in high-elevation regions. ~~For example, over the Alpine region, ERA5 compared to ERA-Interim is characterized by larger daily precipitation amounts during the season December to April, probably in relation to its finer spatial resolution and enhanced orographic processes (average over the period 1980-2014, not shown). ERA5 shows larger precipitation amounts also when compared to~~ provides slightly better performances than GLDAS. The latter has a precipitation bias that is strongly varying from season to season, with large overestimation errors in the observation-based dataset HISTALP as a reference over the same period (not shown). A careful evaluation of these datasets is thus recommended before using them, especially in mountain areas. last two snow seasons (-22%, -25%, -25%, +129%, +102% respectively). By contrast, the ERA5 precipitation bias has smaller fluctuations from season to season (+34%, +16%, +38%, +47%, +41%) resulting in better and more stable performances compared to GLDAS.

The present analysis allows to straightforwardly evaluate the performances of each model with data of gradually lower accuracy. While, as expected ~~, also from previous studies (e.g. Jin et al., 1999; Boone and Etchevers, 2001; Luo et al., 2003; Feng et al., 2008)~~ with accurate forcing the most sophisticated model provides the best agreement with SWE and snow depth observations and the simplest models provide the worst (Fig 6d), more heterogeneous model responses are obtained when lower accuracy data are employed. The most sophisticated model SNOWPACK is not the best performing model throughout all experiments, even though it usually ranks among the best performing ones especially in reproducing snow depth. The simplest snow model considered in the analysis, S3M, is not always the worst model, especially when low accuracy forcings are employed. SMASH shows an interesting behavior, with no brilliant performances with optimal forcing but outperforming many other models when using lower accuracy inputs. ~~Indeed, SMASH ranks among the best performing models in the~~ TIME-12h, MeteoIO, ERA5, ERAI-LR and ERAI-BIAS experiments, suggesting that it can be employed in data-sparse conditions with comparable results as results that are comparable to those of the more sophisticated models.

The GEOTOP model provides the best snow depth estimates when forced by MeteoIO, ERA5 and ERAI-LR. However, all these forcing datasets have a cold temperature bias, and GEOTOP ~~model~~ is affected by a ~~sistematic~~ systematic underestimation

error on snow depth. These errors offset each other, with the effect that the RMSE on snow depth simulations is ~~smallest compared to~~ smaller than for the other models. Conversely, when using ERA-Interim forcing, GEOTOP performances are the worst ones ~~due to~~ owing to the positive temperature bias of the reanalysis dataset, which increases the underestimation of snow depth simulations. In this set of experiments GEOTOP ~~model show~~ shows weaknesses in reproducing ~~both snow-water equivalent and snow density~~ the snow density and depth, thus calling for a check of its snow scheme.

The UTOPIA and HTESSEL models perform as well as the most sophisticated ~~model~~ SNOWPACK with optimal forcing, but they require less input data, for example they do not need ~~surface~~ ground temperature. These models can be employed when no information on snowpack internal structure and stratification is needed. UTOPIA and ~~CHTESSEL~~ HTESSEL provide good performances also with low temporal resolution forcings up to 6 hours and with ERA-Interim forcing. However, lower skills are found when employing ~~other~~ the low-accuracy input dataset (~~TIME-12h, MeteoIO~~) MeteoIO, suggesting that ~~the~~ UTOPIA and HTESSEL ~~models~~ can be sensitive to the bias in the meteorological forcing.

In agreement with former studies (e.g. Essery et al., 2013) also in our analysis the best performing models have i) an explicit representation of the meltwater retention and refreezing in the snowpack and ii) an intermediate-complexity representation of the snow albedo as a function ~~at least of~~ of at least the surface temperature and snow age. According to our results, the representation of the snowpack as a medium with multiple layers alone does not guarantee improved results compared to models with single-layer snow schemes but ~~taking~~ able to take into account meltwater infiltration and refreezing within the snowpack.

This intercomparison exercise has been performed on a single mountain site, Torgnon, ~~in the Western Alps, providing ideal conditions which provides ideal conditions~~ (high-quality input and validation data, low wind speeds) to perform the sensitivity study which we aimed to. Further analysis at other test sites would be useful to explore the extent to which our results could be generalized to different situations or models. We can hypothesize that the effect of the degradation in time of the forcing is probably not site-specific and similar results could be obtained in other sites (see e.g. Ménard et al., 2015). In order to assess the exportability of the results obtained in the reanalysis-driven experiments, in Appendix D we evaluate the biases of the reanalyses considered in this study (ERA5, ERA-Interim and GLDAS) in reproducing the main drivers of the snowpack processes, i.e. temperature and total precipitation, compared to reference datasets (e.g. E-OBS version 13, Haylock et al., 2008) over the entire Greater Alpine Region (GAR, 4°E-19°E, 43°N-49°N). The time-averaged biases found at the Torgnon site are spatially consistent with those found at the mountain range scale, with the magnitude of the bias slightly varying across the region and with elevation. This analysis broadens the perspective beyond the specific case of the Torgnon site and provides a guidance on the exportability of our experiment results to other areas in the Alpine region.

6 Conclusions

Relevant issues in snow modelling are the sparseness of meteorological stations providing all the variables required to drive and validate snow models, and the large uncertainties affecting the available measurements. Moreover, in mountain areas the

spatial variability of the meteorological parameters is high, and ~~the~~ in-situ stations ~~may~~ could be scarcely representative of the conditions in the surrounding areas.

Currently available snow models cover a wide range of complexities, from the most sophisticated schemes that resolve the internal structure of the snowpack to the simplest ones that only provide a coarse estimate of snow depth and snow water equivalent. While several studies evaluate snow models when driven by accurate meteorological data, efforts are still needed to investigate how the models perform when forced by lower-accuracy meteorological data, as are those typically used in mountain areas.

This study evaluates snow models of different complexities assessing their sensitivity to the accuracy of the input data. An interesting result is that some of the simplest models perform equally well or even better than sophisticated models when input data are poor. For example, the intermediate-complexity model SMASH provides lower RMSE values on SWE-snow depth simulations than many other higher-complexity models when driven by 12-hourly data, MeteoIO spatially interpolated data, GLDAS, ERA5, or the bias-adjusted ERA-Interim reanalysis. The lowest-complexity model considered in this study, S3M, provides ~~comparable performances as~~ performances that are comparable to those of the most sophisticated snow model analyzed here, SNOWPACK, when it is driven by bias-adjusted ERA-Interim data.

On the other hand, this study also shows that sophisticated snow models such as SNOWPACK can successfully reproduce snowpack variability across a wider spectrum of conditions compared to simpler snow models, outperforming them in case of isolated snowfall followed by rapid ablation. Sophisticated models provide good and more stable performances across different seasons. It is worth stressing that the most detailed snow model considered here, SNOWPACK, though not being the best performing model throughout all the experiments with lower accuracy forcings, always ranks among the best performing models at reproducing snow depth in all experiments.

Two of the intermediate-complexity snow models, HTESSEL and UTOPIA, ~~provide comparable in the case of optimal forcing~~ provide skills in reproducing SWE and snow depth ~~to that are comparable to those of~~ the most sophisticated model SNOWPACK ~~in case of optimal forcing~~. In addition, they show similar skill across different seasons, thus revealing significant robustness in reproducing a variety of conditions. HTESSEL and UTOPIA can thus be considered a good trade-off between model complexity and model accuracy in case of high-quality forcing data, while they are found to be sensitive to biases in the forcing.

Some properties which are common to all models can be highlighted: i) difficulty in reproducing snow density, especially in late spring at the end of the snow season; ii) low model sensitivity to the use of surrogate radiation input data instead of the measured ones, at least for the test site considered here; iii) comparable performances when driven by 3-hourly or 30-minute data, suggesting the possibility to use lower frequency data (up to 3 hours) without losing accuracy on the snow output; iv) decrease of the models reliability, but not uniformly across the different models, when coarse-grid forcings are employed; v) substantial improvement of the models performances, reducing the differences between models of different complexity, after applying a very simple bias adjustment to temperature (and consistently snowfall) forcing.

The present study ~~exploring the relations between snow model complexity, accuracy of the forcing data and model performance~~ has been conceived to set the basis for high-resolution modeling of mountain snow resources at the catchment and regional

scales in areas where direct meteorological measurements are insufficient or unavailable and one has to rely on coarse resolution forcing. Such sensitivity experiments pave the way for the production of long-term fine-resolution reanalyses for the alpine snowpack, currently identified as a major gap for cryosphere studies (Beniston et al., 2018; Terzago et al., 2017), as well as of high-resolution future projections of the snowpack conditions. In this case snow models can be employed to ~~downscale~~ refine
5 the climate information provided by regional climate models and achieve information on snowpack characteristics at the scales required by hydrological applications, typically below 1 km. This approach, dedicated to the reconstruction of the mountain snowpack variability at fine scales is complementary to the one pursued by the ongoing ESM-SnowMIP initiative (Krinner et al., 2018) aiming at ~~the improvement~~ improving of the representation of snow processes and snow-related climate feedbacks in global climate models. Both approaches address issues which have been highlighted as important in cryospheric sciences
10 (Beniston et al., 2018; Terzago et al., 2017) and provide information for a range of applications including the estimation of climate change impacts on the relevant socio-economic and environmental sectors.

Data availability. The data employed in this study are available upon request.

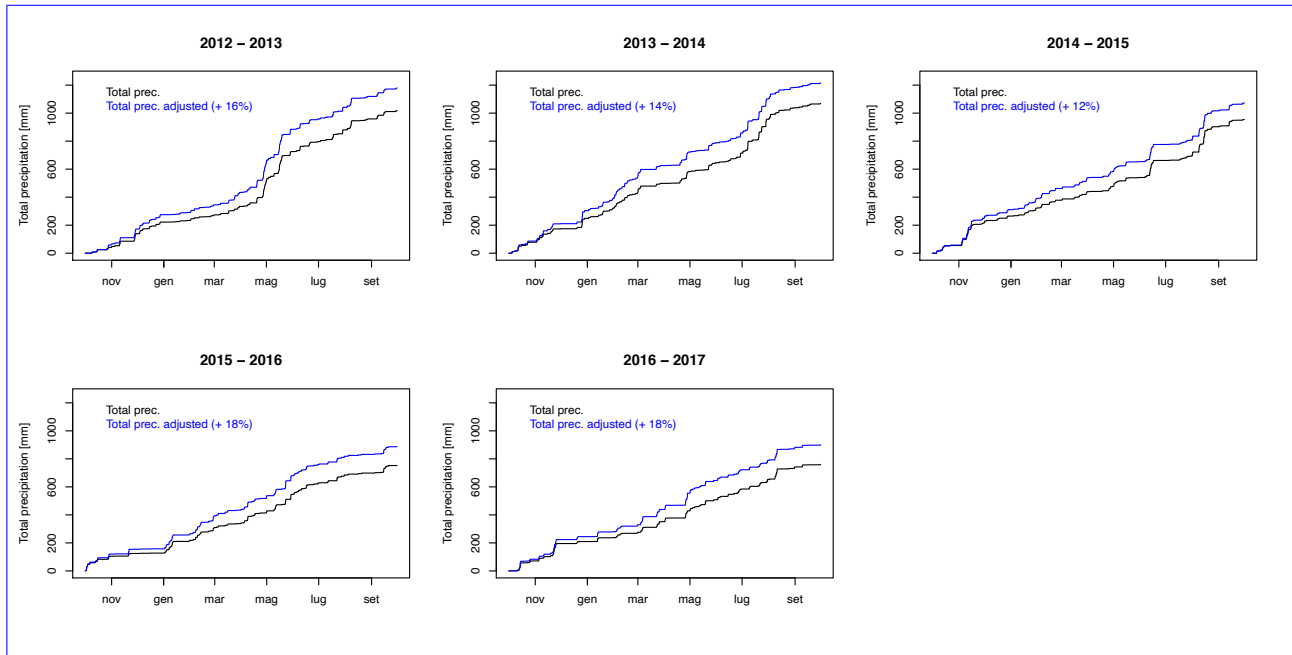


Figure A1. Cumulated total precipitation at the Torgnon site measured by the OTT Pluvio2 precipitation gauge (black) compared to the precipitation adjusted with the Kochendorfer method (blue) for the five snow seasons considered.

Appendix A: **Spatial interpolation of meteorological forcings from neighboring stations** **Uncertainty associated with the precipitation measurements in Torgnon**

We discuss here the uncertainty associated with the observed precipitation and in particular the undercatch of snow which is common in mountain areas. The primary cause for snow precipitation undercatch is related to wind speed, with the amount of precipitation measured by a precipitation gauge relative to the actual amount of precipitation decreasing with increasing wind speed.

We quantified the wind-induced precipitation measurements errors by applying the method described in (Kochendorfer et al., 2017a, b). This method, derived by comparing precipitation measurements from unshielded and shielded (reference) gauges, consists in calculating a catch efficiency (CE), function of air temperature and wind speed, so that the inverse (CE^{-1}) can be used to correct actual precipitation data. The method has been specifically developed for OTT Pluvio2 gauges, i.e. of the same type employed at the Torgnon site.

Figure A1 shows the cumulated total precipitation at the Torgnon site measured by the precipitation gauge (black) compared to the precipitation adjusted with the Kochendorfer method (blue).

The adjusted cumulated total precipitation exceeds the measured precipitation by 16% in average over the 5 snow seasons. As the correction of total precipitation directly affects the amount of solid precipitation, we tested the effects of such correction on snow model simulations. We performed an additional experiment (CTL_prc-adj) in which the model forcing is the same as

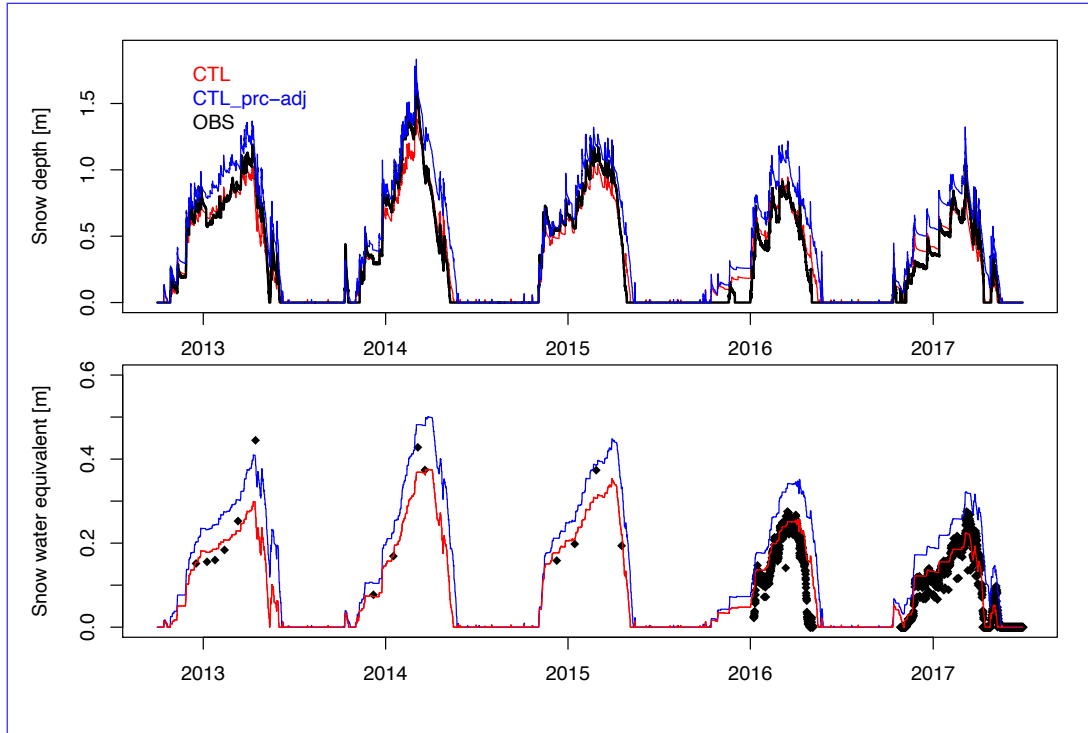


Figure A2. Snow depth (upper panel) and snow water equivalent (lower panel) simulated by the SNOWPACK model when the adjusted total precipitation forcing is employed (CTL_prc-adj) compared to the control run (CTL) and observations.

in the CTL run except for total precipitation, which is now the adjusted one. The snowfall fraction is then calculated from the adjusted total precipitation.

Figure A2 shows the results for the SNOWPACK model, and it displays the simulated snow depth (upper panel) and snow water equivalent (bottom panel) obtained in the CTL and in the CTL_prc-adj runs compared to observations. In all snow seasons the snow depth and the snow water equivalent are remarkably overestimated in the CTL_prc-adj experiment compared both observations and the CTL run. The additional snowfall input derived from the precipitation adjustment leads to an excess of snow accumulation on the ground which can be quantified in an average snow depth bias of 0.17 m compared to the -0.001 m bias in the CTL run. The RMSE is double in the CTL_prc-adj run compared to the CTL run (see Table A1 for details).

As the precipitation adjustment method itself is affected by its own uncertainties, and given that the application of the precipitation adjustment leads to a worsening in the snow model performances, we decide to employ the original precipitation measurements as forcing in the snow model experiments.

Table A1. [SNOWPACK model RMSE and bias for the simulated snow depth and snow water equivalent variables in the CTL_prc-adj experiment and in the control run \(CTL\).](#)

	Snow depth		SWE	
	RMSE [m]	BIAS [m]	RMSE [m]	BIAS [m]
CTL	0.10	-0.001	0.04	0.02
CTL_prc-adj	0.20	0.170	0.10	0.09

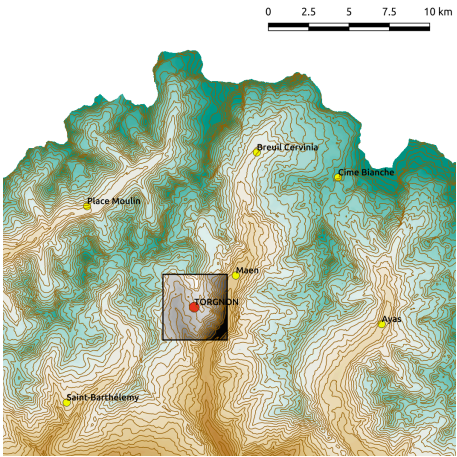


Figure B1. Location of the 6 neighboring stations used for producing the interpolated dataset for the MeteoIO experiment. The grey square represents the extension of the digital elevation model used for the interpolation.

Appendix B: [Spatial interpolation of meteorological forcings from neighboring stations](#)

In hydrological and snow modeling the spatial interpolation of ground meteorological observations is commonly employed to derive spatially continuous meteorological forcing to drive the models. In this work, we evaluate the response of snow models with such forcing. An interpolated dataset for Torgnon monitoring site has been prepared exploiting the MeteoIO library (Bavay and Egger, 2014). The meteorological data are interpolated from six neighboring stations, over a squared digital elevation model of 16 km² with a grid resolution of 50 meters centered on the coordinate of [the](#) Torgnon monitoring site (Fig. B1 and Tab. [A1B1](#)).

Table B1. Characteristics of the meteorological stations used for the spatial interpolation with MeteoIO library and measured parameters: TA = air temperature; PTOT = precipitation (OTT); SWIN = incoming short wave solar radiation; VW-DW = wind speed and direction; RH = relative humidity. The stations belong to the regional meteorological network of the Aosta Valley.

Station name	Elevation [m a.s.l.]	Distance [km]	TA	PTOT	SWIN	VW-DW	RH
Cime Bianche	3100	12	x	x	x	x	x
Saint-Berthélemy	1675	9.8	x		x	x	x
Place Moulin	1980	9.1	x	x	x	x	x
Breuil Cervinia	2000	10.3	x			x	x
Maen	1310	3.2	x			x	x
Ayas	1566	11.6	x			x	x

Appendix C: [The impact of the time interpolation method for SWIN in the TIME-12h experiment](#)

We test the impact of using two different methods to derive 30 minute temporal resolution shortwave incoming radiation input when only measurements at 00:00 and 12:00 UTC+01:00 are available (as in the TIME-12h experiment). The first method is a basic linear interpolation of the available measurements. The second method is slightly more sophisticated and employs the potential (clear-sky) incoming shortwave radiation (Knauer et al., 2018) at 30 minute temporal resolution and at the coordinates of the Torgnon station, and the SWIN station measurements at 12:00. For each day of the year, the 48 daily values of potential radiation are rescaled according to the observed SWIN value at 12:00, to obtain an “estimated SWIN” (Figure C1a).

With the first method, based on the linear interpolation, the average difference between the estimated and the observed SWIN radiation over the full period is large ($+97 \text{ W/m}^2$) while with the second method, based on the scaling of the potential radiation, the difference is close to zero (-0.87 W/m^2).

In order to test the impact of the method to interpolate SWIN radiation on snow simulations, we run two experiments in which the forcing is the Torgnon data sampled every 12 hours as explained in Section 4. The two forcings differ for the SWIN radiation input: in one case obtained by linearly interpolating SWIN measurements (TIME-12h-LIN) and in the other case obtained by rescaling the potential radiation as explained above (TIME-12h-SWINPOT).

Figure C1b shows the results of the two experiments, TIME-12h-LIN and TIME-12h-SWINPOT, compared to the CTL run and observations, for the SNOWPACK model and for the snow depth variable. The use of the SWIN forcing derived from the potential radiation leads to a remarkable improvement in the agreement with observations compared to the case when linearly interpolated SWIN is used, with the model RMSE reduced to a value which is comparable to that obtained in the CTL run (Table C1). The results for all snow models are reported in Table C1.

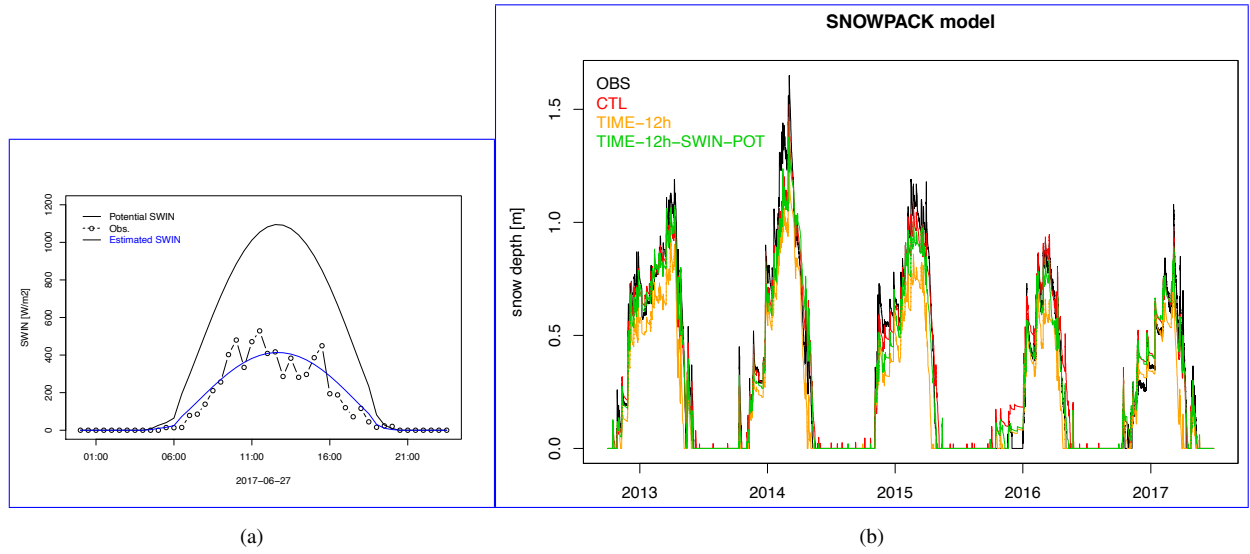


Figure C1. a) Measured shortwave incoming radiation (SWIN) at the Torgnon site for the day 26 June 2017 (points), potential SWIN for that day (solid black line), “estimated SWIN” from the scaling of the potential SWIN on the value registered at h 12:00; b) Snow depth simulations obtained with the SNOWPACK model for the experiment TIME-12h-SWIN-POT compared to TIME-12h-LIN, the CTL run and observations.

Table C1. Model RMSE for the simulated snow depth in the CTL run, the TIME-12h-LIN and the TIME-12h-SWIN-POT experiments, compared to observations.

Model	RMSE snow depth [m]		
	CTL	TIME-12h-LIN	TIME-12h-SWIN-POT
SNOWPACK	0.10	0.21	0.11
GEOTOP	0.17	0.37	0.35
HTESSEL	0.11	0.44	0.12
UTOPIA	0.12	0.38	0.26
SMASH	0.17	0.17	0.17
S3M	0.25	0.38	0.39

Appendix D: Exportability of the results from reanalysis-driven experiments

In order to address the issue of the exportability of the methods and results of the reanalysis-driven experiments to other areas of the Alps, we evaluated the biases of the reanalyses in reproducing the main drivers of the snowpack processes, i.e. temperature

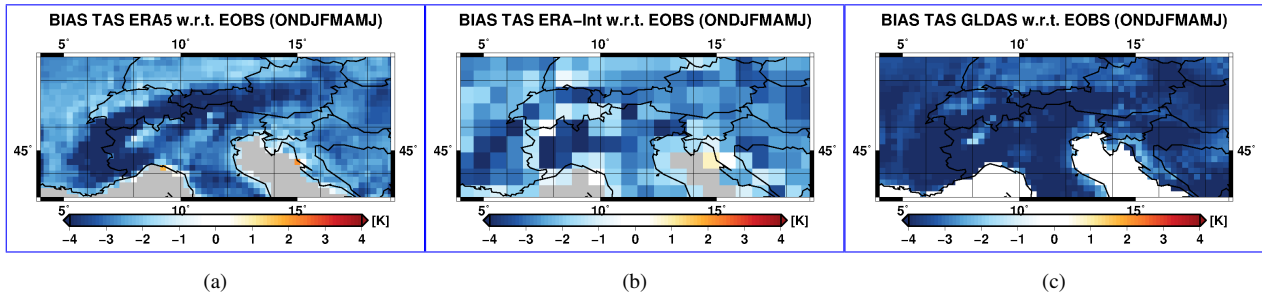


Figure D1. BIAS of ERA5, ERA-Interim and GLDAS air temperatures with respect to E-OBS observations over the Greater Alpine Region. Temperatures have been averaged over the months from October to June and over the period 1980-2014 in the case of ERA5 and ERA-Interim, over the period 2000-2014 in the case of GLDAS.

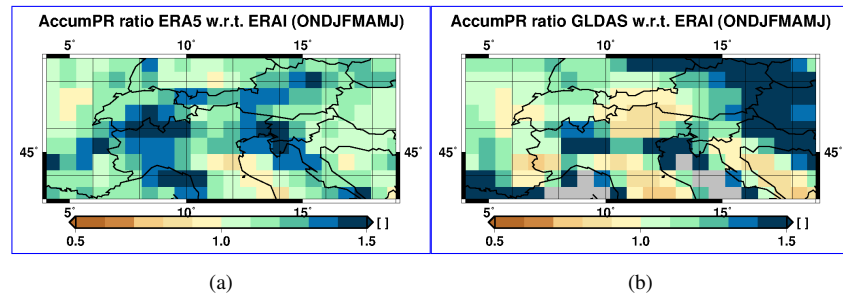


Figure D2. ERA5 and GLDAS relative differences with respect to ERA-Interim for the October-June accumulated precipitation over the periods 1980-2014 and 2000-2014 respectively.

and total precipitation, compared to observational data. The aim is to evaluate the spatial distribution of the temperature and precipitation biases and their consistency at the mountain range scale.

ERA5, ERA-Interim and GLDAS temperatures have been averaged over the months October-June and over the years 1980-2014 (except for GLDAS which is available since 2000 only, so the averages have been calculated over the period 2000-2014), and then compared to the observational dataset E-OBS version 13 (Haylock et al., 2008) over the Greater Alpine Region (GAR, 4°E-19°E, 43°N-49°N). E-OBS is a daily gridded data set at 0.25° resolution, based on the European Climate Assessment and Data set station measurements.

ERA5 and GLDAS temperature biases are large and negative over the entire GAR (Figure D1). GLDAS bias is especially strong and it exceeds -4°C in most of the region. ERA5 bias is large at high elevations and decreases towards the lowlands. Compared to ERA5 and GLDAS, ERA-Interim temperature is in better agreement with observations, with mainly negative biases across the region and values close to zero (both positive and negative values) at the mountain ridges in Western Alps. All these results are consistent with those found at the Torgnon site (Fig. 3), so the biases at the point scale are reflected at the mountain range scale.

Regarding precipitation, it is well known that standard surface station gauges have problems in capturing snowfall and thus they underestimate total precipitation in mountain areas. Similarly, also observational-based dataset such as E-OBS have been found to suffer the underestimation of precipitation at high elevations (Turco et al., 2013). To overcome this problem, instead of using observation-based datasets as a reference, we evaluate precipitation ratios with respect to a reanalysis (ERA-Interim), which inherently takes into account orographic effects. Figure D2 shows the ERA5 and GLDAS October-to-June accumulated precipitation ratios relative to ERA-Interim over the periods 1980-2014 and 2000-2014 respectively (GLDAS is available since 2000). Also in this case ERA5 spatial pattern is homogeneous over the Alpine range, with ERA5 showing consistently more precipitation than ERA-Interim in the mountain areas. GLDAS precipitation is found to be in slightly better agreement with the ERA-Interim reanalysis than ERA5, with relative precipitation bias close to 1 over the Alpine range.

Overall, this analysis providing information on the spatial variability of the temperature and precipitation biases in the reanalysis products over the Alpine region broadens the perspective beyond the specific case of the Torgnon site. The time-averaged biases at the Torgnon site result spatially consistent with those found at the mountain range scale, with the magnitude of the bias slightly varying across the region and with the elevation. Similar biases in the forcing suggest that the methods applied in the reanalysis-driven experiments could be extended to other sites in the Alps and could lead to results not too dissimilar from those found at Torgnon.

Author contributions. ST, AP, CC, EC, SG, UMC, PP conceived the idea of the experiments. All authors participated in the collection of the meteorological datasets for the experiments. ST, VA, GA, LC, DD, GP, PP performed the simulations. ST analyzed the simulations and prepared all figures, all authors provided support in the interpretation of the results. ST wrote the paper with support from all authors.

Competing interests. The authors declare that no competing interests are present.

20 *Disclaimer.* TEXT

Acknowledgements. This work received funding from the Italian Project of Interest NextData of the Italian Ministry for Education, University and Research and from the European Union's Horizon 2020 research and innovation program under grant agreement no. 641762 (ECOPOTENTIAL). Part of this work was performed in the framework of the MEDSCOPE (MEDiterranean Services Chain based On climate PrEdictions) ERA4CS project (grant agreement no. 690462) funded by the European Union. ~~ERA5: Copernicus Climate Change Service (C3S) (2017): ERA5: Fifth generation of ECMWF atmospheric reanalyses of the global climate . Copernicus Climate Change Service Climate Data Store (CDS), date of citation.~~ Jost von Hardenberg acknowledges support from the European Union's Horizon 2020 research and innovation programme under Grant agreement 641816 (CRESCENDO).

References

- Balsamo, G., Beljaars, A., Scipal, K., Viterbo, P., van den Hurk, B., Hirschi, M., and Betts, A. K.: A revised hydrology for the ECMWF model: Verification from field site to terrestrial water storage and impact in the Integrated Forecast System, *Journal of Hydrometeorology*, 10, 623–643, <https://doi.org/10.1175/2008JHM1068.1>, 2009.
- 5 Bartelt, P. and Lehning, M.: A physical SNOWPACK model for the Swiss avalanche warning: Part I: numerical model, *Cold Regions Science and Technology*, 35, 123–145, 2002.
- Bavay, M. and Egger, T.: MeteoIO 2.4. 2: a preprocessing library for meteorological data, *Geoscientific Model Development*, 7, 3135–3151, 2014.
- Beniston, M., Farinotti, D., Stoffel, M., Andreassen, L. M., Coppola, E., Eckert, N., Fantini, A., Giacona, F., Hauck, C., Huss, M., Huwald, H., Lehning, M., López-Moreno, J.-I., Magnusson, J., Marty, C., Morán-Tejeda, E., Morin, S., Naaim, M., Provenzale, A., Rabatel, A., Six, D., Stötter, J., Strasser, U., Terzago, S., and Vincent, C.: The European mountain cryosphere: a review of its current state, trends, and future challenges, *Cryosphere*, 12, 759–794, 2018.
- 10 Boni, G., Castelli, F., Gabellani, S., Machiavello, G., and Rudari, R.: Assimilation of MODIS snow cover and real time snow depth point data in a snow dynamic model, in: *Geoscience and Remote Sensing Symposium (IGARSS), 2010 IEEE International*, pp. 1788–1791, IEEE, 2010.
- 15 Boone, A. and Etchevers, P.: An intercomparison of three snow schemes of varying complexity coupled to the same land surface model: Local-scale evaluation at an Alpine site, *Journal of Hydrometeorology*, 2, 374–394, 2001.
- Boone, A., Habets, F., Noilhan, J., Clark, D., Dirmeyer, P., Fox, S., Gusev, Y., Haddeland, I., Koster, R., Lohmann, D., et al.: The Rhone-aggregation land surface scheme intercomparison project: An overview, *Journal of Climate*, 17, 187–208, 2004.
- 20 Bowling, L. C., Lettenmaier, D. P., Nijssen, B., Graham, L., Clark, D. B., El Maayar, M., Essery, R., Goers, S., Gusev, Y. M., Habets, F., et al.: Simulation of high-latitude hydrological processes in the Torne–Kalix basin: PILPS Phase 2 (e): 1: Experiment description and summary intercomparisons, *Global and Planetary Change*, 38, 1–30, 2003.
- Cassardo, C.: UTOPIA: The Manual of Version 2015, Università di Torino, <https://doi.org/10.13140/RG.2.2.29664.38404>, https://www.researchgate.net/publication/323200198_UTOPIA_The_Manual_of_Version_2015, 2015.
- 25 Clark, M. P., Kavetski, D., and Fenicia, F.: Pursuing the method of multiple working hypotheses for hydrological modeling, *Water Resources Research*, 47, 2011.
- Dee, D. P., Uppala, S. M., Simmons, A. J., Berrisford, P., Poli, P., Kobayashi, S., Andrae, U., Balmaseda, M. A., Balsamo, G., Bauer, P., Bechtold, P., Beljaars, A. C. M., van de Berg, L., Bidlot, J., Bormann, N., Delsol, C., Dragani, R., Fuentes, M., Geer, A. J., Haimberger, L., Healy, S. B., Hersbach, H., Hólm, E. V., Isaksen, L., Kållberg, P., Köhler, M., Matricardi, M., McNally, A. P., Monge-Sanz, B. M., Morcrette, J.-J., Park, B.-K., Peubey, C., de Rosnay, P., Tavolato, C., Thépaut, J.-N., and Vitart, F.: The ERA-Interim reanalysis: configuration and performance of the data assimilation system, *Quarterly Journal of the Royal Meteorological Society*, 137, 553–597, <https://doi.org/10.1002/qj.828>, <http://dx.doi.org/10.1002/qj.828>, 2011.
- 30 Di Mauro, B., Garzonio, R., Rossini, M., Filippa, G., Pogliotti, P., Galvagno, M., Morra di Cella, U., Migliavacca, M., Baccolo, G., Clemenza, M., Delmonte, B., Maggi, V., Dumont, M., Tuzet, F., Lafaysse, M., Morin, S., Cremonese, E., and Colombo, R.: Saharan dust events in the European Alps: role on snowmelt and geochemical characterization, *The Cryosphere Discussions*, 2018, 1–28, <https://doi.org/10.5194/tc-2018-241>, <https://www.the-cryosphere-discuss.net/tc-2018-241/>, 2018.
- 35

- Dutra, E., Balsamo, G., Viterbo, P., Miranda, P. M., Beljaars, A., Schär, C., and Elder, K.: An improved snow scheme for the ECMWF land surface model: description and offline validation, *J. Hydrometeor.*, 11, 899–916, <https://doi.org/10.1175/2010jhm1249.1>, 2010.
- Dutra, E., Viterbo, P., Miranda, P. M., and Balsamo, G.: Complexity of snow schemes in a climate model and its impact on surface energy and hydrology, *Journal of Hydrometeorology*, 13, 521–538, 2012.
- 5 Endrizzi, S., Gruber, S., Dall’Amico, M., and Rigon, R.: GEOTop 2.0: simulating the combined energy and water balance at and below the land surface accounting for soil freezing, snow cover and terrain effects, *Geoscientific Model Development*, 7, 2831–2857, 2014.
- Essery, R., Morin, S., Lejeune, Y., and Ménard, C. B.: A comparison of 1701 snow models using observations from an alpine site, *Advances in Water Resources*, 55, 131 – 148, <https://doi.org/http://dx.doi.org/10.1016/j.advwatres.2012.07.013>, <http://www.sciencedirect.com/science/article/pii/S0309170812002011>, snow–Atmosphere Interactions and Hydrological Consequences, 2013.
- 10 Etchevers, P., Martin, E., Brown, R., Fierz, C., Lejeune, Y., Bazile, E., Boone, A., Dai, Y., Essery, R., Fernandez, A., et al.: SnowMIP, an intercomparison of snow models: first results, in: *Proceedings of the International Snow Science Workshop*, Penticton, Canada, vol. 29, 2002.
- Etchevers, P., Martin, E., Brown, R., Fierz, C., Lejeune, Y., Bazile, E., Boone, A., Dai, Y.-J., Essery, R., Fernandez, A., et al.: Validation of the energy budget of an alpine snowpack simulated by several snow models (SnowMIP project), *Annals of Glaciology*, 38, 150–158, 15 2004.
- Feng, X., Sahoo, A., Arsenault, K., Houser, P., Luo, Y., and Troy, T. J.: The impact of snow model complexity at three CLPX sites, *Journal of Hydrometeorology*, 9, 1464–1481, 2008.
- Filippa, G., Cremonese, E., Galvagno, M., Migliavacca, M., Di Cella, U. M., Petey, M., and Siniscalco, C.: Five years of phenological monitoring in a mountain grassland: inter-annual patterns and evaluation of the sampling protocol, *International journal of biometeorology*, 20 59, 1927–1937, 2015.
- Galvagno, M., Wohlfahrt, G., Cremonese, E., Rossini, M., Colombo, R., Filippa, G., Julitta, T., Manca, G., Siniscalco, C., Morra di Cella, U., et al.: Phenology and carbon dioxide source/sink strength of a subalpine grassland in response to an exceptionally short snow season, *Environmental Research Letters*, 8, 025 008, 2013.
- Haylock, M., Hofstra, N., Klein Tank, A., Klok, E., Jones, P., and New, M.: A European daily high-resolution gridded data set of surface 25 temperature and precipitation for 1950–2006, *Journal of Geophysical Research: Atmospheres*, 113, 2008.
- Hersbach, H. and Dee, D.: ERA5 reanalysis is in production, *ECMWF Newsletter* 147, ECMWF, 2016.
- Jin, J., Gao, X., Yang, Z.-L., Bales, R., Sorooshian, S., Dickinson, R. E., Sun, S., and Wu, G.: Comparative analyses of physically based snowmelt models for climate simulations, *Journal of Climate*, 12, 2643–2657, 1999.
- Knauer, J., El-Madany, T. S., Zaehle, S., and Migliavacca, M.: Bigleaf—An R package for the calculation of physical and physiological 30 ecosystem properties from eddy covariance data, *PloS one*, 13, 2018.
- Kochendorfer, J., Nitu, R., Wolff, M., Mekis, E., Rasmussen, R., Baker, B., Earle, M. E., Reverdin, A., Wong, K., Smith, C. D., Yang, D., Roulet, Y.-A., Buisan, S., Laine, T., Lee, G., Aceituno, J. L. C., Alastrué, J., Isaksen, K., Meyers, T., Brækkan, R., Landolt, S., Jachcik, A., and Poikonen, A.: Analysis of single-Alter-shielded and unshielded measurements of mixed and solid precipitation from WMO-SPICE, *Hydrology and Earth System Sciences*, 21, 3525–3542, <https://doi.org/10.5194/hess-21-3525-2017>, <https://www.hydrol-earth-syst-sci.net/21/3525/2017/>, 2017a.
- 35 Kochendorfer, J., Rasmussen, R., Wolff, M., Baker, B., Hall, M. E., Meyers, T., Landolt, S., Jachcik, A., Isaksen, K., Brækkan, R., et al.: The quantification and correction of wind-induced precipitation measurement errors, *Hydrology and Earth System Sciences*, 21, 1973–1989, 2017b.

- Krinner, G., Derksen, C., Essery, R., Flanner, M., Hagemann, S., Clark, M., Hall, A., Rott, H., Brutel-Vuilmet, C., Kim, H., Ménard, C. B., Mudryk, L., Thackeray, C., Wang, L., Arduini, G., Balsamo, G., Bartlett, P., Boike, J., Boone, A., Chérut, F., Colin, J., Cuntz, M., Dai, Y., Decharme, B., Derry, J., Ducharne, A., Dutra, E., Fang, X., Fierz, C., Ghattas, J., Gusev, Y., Haverd, V., Kontu, A., Lafaysse, M., Law, R., Lawrence, D., Li, W., Marke, T., Marks, D., Nasonova, O., Nitta, T., Niwano, M., Pomeroy, J., Raleigh, M. S., Schaedler, G., Semenov, V., Smirnova, T., Stacke, T., Strasser, U., Svenson, S., Turkov, D., Wang, T., Wever, N., Yuan, H., and Zhou, W.: ESM-SnowMIP: Assessing models and quantifying snow-related climate feedbacks, *Geoscientific Model Development Discussions*, 2018, 1–32, <https://doi.org/10.5194/gmd-2018-153>, <https://www.geosci-model-dev-discuss.net/gmd-2018-153/>, 2018.
- Kumar, M., Marks, D., Dozier, J., Reba, M., and Winstral, A.: Evaluation of distributed hydrologic impacts of temperature-index and energy-based snow models, *Advances in Water Resources*, 56, 77–89, 2013.
- 10 Lawrence, M. G.: The Relationship between Relative Humidity and the Dewpoint Temperature in Moist Air: A Simple Conversion and Applications, *Bulletin of the American Meteorological Society*, 86, 225–233, <https://doi.org/10.1175/bams-86-2-225>, <http://dx.doi.org/10.1175/bams-86-2-225>, 2005.
- Luo, L., Robock, A., Vinnikov, K. Y., Schlosser, C. A., Slater, A. G., Boone, A., Etchevers, P., Habets, F., Noilhan, J., Braden, H., et al.: Effects of frozen soil on soil temperature, spring infiltration, and runoff: Results from the PILPS 2 (d) experiment at Valdai, Russia, *Journal of Hydrometeorology*, 4, 334–351, 2003.
- 15 Magnusson, J., Wever, N., Essery, R., Helbig, N., Winstral, A., and Jonas, T.: Evaluating snow models with varying process representations for hydrological applications, *Water Resources Research*, 51, 2707–2723, 2015.
- Ménard, C. B., Ikonen, J., Rautiainen, K., Aurela, M., Arslan, A. N., and Pulliainen, J.: Effects of meteorological and ancillary data, temporal averaging, and evaluation methods on model performance and uncertainty in a land surface model, *Journal of Hydrometeorology*, 16, 2559–2576, 2015.
- 20 Piazz, G., Thirel, G., Campo, L., and Gabellani, S.: A particle filter scheme for multivariate data assimilation into a point-scale snowpack model in an Alpine environment, *The Cryosphere*, 12, 2287–2306, 2018.
- Piazz, G., Campo, L., Gabellani, S., Castelli, F., Cremonese, E., di Cella, U. M., Stevenin, H., and Ratto, S. M.: An EnKF-based scheme for snow multivariable data assimilation at an Alpine site, *Journal of Hydrology and Hydromechanics*, 67, 4–19, 2019.
- 25 Raleigh, M., Lundquist, J., and Clark, M.: Exploring the impact of forcing error characteristics on physically based snow simulations within a global sensitivity analysis framework, *Hydrology and Earth System Sciences*, 19, 3153–3179, 2015.
- Rodell, M., Houser, P., Jambor, U., Gottschalck, J., Mitchell, K., Meng, C.-J., Arsenault, K., Cosgrove, B., Radakovich, J., Bosilovich, M., et al.: The global land data assimilation system, *Bulletin of the American Meteorological Society*, 85, 381–394, 2004.
- Rui, H. and Beaudoin, H.: README Document for NASA GLDAS Version 2 Data Products, Goddard Earth Sciences Data and Information Services Center (GES DISC), Greenbelt, USA, 2018.
- 30 Rutter, N., Essery, R., Pomeroy, J., Altimir, N., Andreadis, K., Baker, I., Barr, A., Bartlett, P., Boone, A., Deng, H., Douville, H., Dutra, E., Elder, K., Ellis, C., Feng, X., Gelfan, A., Goodbody, A., Gusev, Y., Gustafsson, D., Hellström, R., Hirabayashi, Y., Hirota, T., Jonas, T., Koren, V., Kuragina, A., Lettenmaier, D., Li, W.-P., Luce, C., Martin, E., Nasonova, O., Pumpanen, J., Pyles, R. D., Samuelsson, P., Sandells, M., Schädler, G., Shmakin, A., Smirnova, T. G., Stähli, M., Stöckli, R., Strasser, U., Su, H., Suzuki, K., Takata, K., Tanaka, K., Thompson, E., Vesala, T., Viterbo, P., Wiltshire, A., Xia, K., Xue, Y., and Yamazaki, T.: Evaluation of forest snow processes models (SnowMIP2), *Journal of Geophysical Research: Atmospheres*, 114, n/a–n/a, <https://doi.org/10.1029/2008JD011063>, <http://dx.doi.org/10.1029/2008JD011063>, 2009.
- 35

- Schlosser, C. A., Slater, A. G., Robock, A., Pitman, A. J., Vinnikov, K. Y., Henderson-Sellers, A., Speranskaya, N. A., and Mitchell, K.: Simulations of a boreal grassland hydrology at Valdai, Russia: PILPS Phase 2 (d), *Monthly Weather Review*, 128, 301–321, 2000.
- Sims, E. M. and Liu, G.: A parameterization of the probability of snow–rain transition, *Journal of Hydrometeorology*, 16, 1466–1477, 2015.
- Slater, A. G., Schlosser, C. A., Desborough, C., Pitman, A., Henderson-Sellers, A., Robock, A., Vinnikov, K. Y., Entin, J., Mitchell, K., Chen, F., et al.: The representation of snow in land surface schemes: Results from PILPS 2 (d), *Journal of Hydrometeorology*, 2, 7–25, 2001.
- 5 Taylor, K. E.: Summarizing multiple aspects of model performance in a single diagram, *Journal of Geophysical Research: Atmospheres*, 106, 7183–7192, 2001.
- Terzago, S., von Hardenberg, J., Palazzi, E., and Provenzale, A.: Snow water equivalent in the Alps as seen by gridded data sets, CMIP5 and CORDEX climate models, *The Cryosphere*, 11, 1625–1645, <https://doi.org/10.5194/tc-11-1625-2017>, <https://www.the-cryosphere.net/11/1625/2017/>, 2017.
- 10 Turco, M., Zollo, A., Ronchi, C., De Luigi, C., and Mercogliano, P.: Assessing gridded observations for daily precipitation extremes in the Alps with a focus on northwest Italy., *Natural Hazards & Earth System Sciences*, 13, 2013.
- Vionnet, V., Brun, E., Morin, S., Boone, A., Faroux, S., Le Moigne, P., Martin, E., and Willemet, J.: The detailed snowpack scheme Crocus and its implementation in SURFEX v7. 2, *Geoscientific Model Development*, 5, 773–791, 2012.
- 15 Yang, K., Huang, G., and Tamai, N.: A hybrid model for estimating global solar radiation, *Solar energy*, 70, 13–22, 2001.
- Yang, K., Koike, T., and Ye, B.: Improving estimation of hourly, daily, and monthly solar radiation by importing global data sets, *Agricultural and Forest Meteorology*, 137, 43–55, 2006.