# On the reliability of seasonal snow forecasts

Ekaterina Vorobeva[1], Yvan Orsolini[1], Patricia de Rosnay[2], Jonathan Day[2], Retish Senan[2],
Damien Decremer[2], and Frederic Vitart[2]

[1]NILU, Kjeller, Norway
[2]European Centre for Medium-Range Weather Forecasts, Reading, UK

**Correspondence:** Ekaterina Vorobeva (evor@nilu.no)

**Abstract.** Reliable information on seasonal snow conditions is important for long-range weather forecasting and climate modeling. The reliability of winter-mean hindcasts of snow water equivalent (SWE) produced by the ECMWF for the period 1993 – 2022 within the CopERnIcus climate change Service Evolution (CERISE) project is evaluated in this study. In probabilistic forecasting, reliability for a binary event is defined as the consistency between forecast probabilities and observed frequencies. Here, reliability is assessed using two independent SWE datasets (ERA5-Land and ESA Snow-CCI v4) across eight land regions in the Northern Hemisphere excluding mountainous regions. The reliability assessment is performed for two tercile-based binary events representing low- and high snow accumulation winters. Reliability is quantified using a weighted linear regression applied to reliability diagrams and is grouped into five categories from perfect to dangerous. The results show good reliability of the ECMWF seasonal snow hindcasts for both low- and high-snow conditions. The assessment shows sensitivity to the choice of verification dataset, with ERA5-Land yielding slightly higher reliability categories than ESA Snow-CCI. It is found that differences in hindcasts reliability between regions and between verification datasets may be linked to snow variability, model representation, and observational uncertainty.

## 1 Introduction

Snow plays a crucial role for the Earth's surface energy budget of the hydrological and climate systems, making it an important component in climate simulations and long-range weather forecasts. The amount of water stored in a snowpack, known as snow water equivalent (SWE), is a key parameter for flood risk assessment, water resource management, and, more broadly, land-atmosphere coupling and climate modeling. In the warming world, snow conditions are undergoing substantial changes. The global snow cover extent has been shown to decline over the last decades in both models and observations (Bormann et al., 2018; Mudryk et al., 2020; Fox-Kemper et al., 2021) and the length of the snow season has been shown to diminish in mountain regions (Notarnicola, 2020). Snow droughts, either driven by deficits in snowy precipitation or by premature snowmelts, have been linked to hydrological extremes (Zhang et al., 2025). On local scales, a reduction of the fraction of solid to liquid precipitation in mountain regions has also been linked to rainfall extremes and associated hazards (Ombadi et al., 2023).

Given the changes in snow accumulation, melt, and associated hydrological hazards, reliable seasonal predictions of snow conditions are becoming increasingly important. Seasonal snow forecasts are routinely provided by operational meteorological

prediction centres, using state-of-the-art ensemble dynamical prediction systems, based on coupled atmosphere-ocean models. In such forecasts, snow variables are initialized as realistically as possible through data assimilation (DA) methods of varying complexity, which ingest in-situ and space-borne observations. There are currently considerable efforts to improve the fidelity and reliability of such seasonal forecasts, as well as the initialization of land variables, including snow.

30      Realistic initial snow conditions can improve forecast skill not only for snow itself but also for other surface and atmospheric variables through its high albedo and low thermal conductivity, as well as its delayed influence on soil moisture (Orsolini et al., 2013; Jeong et al., 2013; Li et al., 2019). Improved forecasts of snow variables offer hence the potential to improve the predicted circulation, particularly in regions of effective snow-atmosphere coupling, like East Asia (Komatsu et al., 2023). However, seasonal forecasting models have limitations, such as various parametrizations or their relatively coarse horizontal

35 resolution that does not allow resolving snow heterogeneities. Similarly, the data ingested through DA has limitations, such as the sparsity and non-representativeness of in-situ station data, or the observational constraints and errors of remotely sensed data. The DA approaches used in prediction centres have also limitations e.g., assimilation might not be performed in mountain regions.

     Here, we try to address the following question: How reliable are seasonal forecasts of snow? In probabilistic forecasting,

40 reliability refers to the agreement between the predicted forecast probabilities and the corresponding observed frequencies of a binary event (i.e., an event that either occurs or does not occur). Reliability of seasonal-mean near-surface temperature and precipitation forecasted by the European Centre for Medium-Range Weather Forecasts (ECMWF) Integrated Forecasting System (IFS) 4 was examined in Weisheimer and Palmer (2014). However, such assessment have never been performed for probabilistic snow forecasts. In the context of changing snow conditions, reliable forecasting of some aspects of snow phenology (e.g.

45 snow accumulation during cold season, or ablation then after) would ultimately be critical for hydrological applications.

     The estimated reliability of probabilistic snow forecasts will, naturally, depend on the snow product they are assessed against, e.g., outputs from coupled reanalysis systems, standalone land model simulations driven by meteorological forcings, satellite-derived snow estimates or actual snow observations. In the former case, an issue might be that the reanalysis used as verification dataset has the same origin as the one used to initialize the forecast. Numerous snow products are now available to meet the

50 growing demand for historical snow datasets in climate, hydrology and cryosphere research. Recent study by Mudryk et al. (2025) compared twenty three historical gridded snow products against a combination of in-situ snow course and airborne gamma measurements. Based on overall assessment, the best performing product was shown to be the ERA5-Land reanalysis (Muñoz Sabater et al., 2021) albeit it is an off-line surface model run which does not assimilate snow measurements. However, for non-mountainous Eurasia, the European Space Agency (ESA) Climate Change Initiative (CCI) Snow project

55 (hereafter ESA Snow-CCI) product, combining space-born observations of passive microwave brightness temperatures and in-situ measurements of snow depth, scored second-best.

     In this study, the reliability assessment is applied to seasonal snow hindcasts produced by the ECMWF in the framework of the CopERnIcus climate change Service Evolution (CERISE) project (see Acknowledgments and https://cerise-project.eu/), which is aimed at enhancing the quality of the Copernicus Climate Change Service (C3S) reanalysis and seasonal forecast

60 products by means of improved land-atmosphere data assimilation and land surface initialization. The winter-mean hindcasts

of snow water equivalent are evaluated against two products, namely, the ERA5-Land reanalysis and ESA Snow-CCI version 4. The reliability assessment is performed over the years 1993 – 2022 in eight land regions in the Northern Hemisphere excluding mountainous regions.

The paper is organized as following. An introduction into the basics of the reliability assessment is given in Sect. 2.1. All
65   datasets used for the analysis are described in Sect. 2.2–2.4. The results are presented in Sect. 3 and discussed in Sect.4. The conclusions are drawn in Sect. 5.

## 2   Data and Methods

### 2.1   What is a reliable forecast?

The Brier score is a simple measure of error in probabilistic forecasting (Brier, 1950). It is often used in binary situations when
70   an event of interest either occurs or not, for example snow amount exceeds 1 cm of SWE in a particular geographical location. It calculates the mean squared error between the forecast probabilities and observed outcomes and is expressed as

$$\text{BS} = \frac{1}{n} \sum_{i=1}^{n} (f_i - o_i)^2 \tag{1}$$

where $n$ is the number of forecasts, $f_i$ is the forecast probability of the considered binary event for the $i$th forecast, and $o_i$ is the associated outcome (1 if event occurs in verification data and 0 if it does not). One can decompose the Brier score into the
75   uncertainty (UNC), reliability (REL) and resolution (RES) components as

$$\text{BS} = \bar{o}(1 - \bar{o}) + \frac{1}{n} \sum_{k=1}^{m} n_k (f_k - \bar{o}_k)^2 - \frac{1}{n} \sum_{k=1}^{m} n_k (\bar{o}_k - \bar{o})^2 = \text{UNC} + \text{REL} - \text{RES} \tag{2}$$

where $\bar{o}$ is the climatological probability of the event, $m$ is the number of probability bins into which the forecasts are grouped, $f_k$ is the average forecast probability in bin $k$, and $\bar{o}_k$ is the relative frequency of occurrence in bin $k$. From Eq. 2, it is clear that only REL and RES terms reflect the forecast performance.

80   The level of the Brier score improvement compared to that of a reference forecast is represented via the Brier skill score

$$\text{BSS} = \frac{\text{BS}_{ref} - \text{BS}}{\text{BS}_{ref}} \tag{3}$$

here, positive BSS values indicate that the forecast performs better than the reference forecast. It is a common practice to use the forecast climatology ($f_k = \bar{o}_k = \bar{o}$) as the baseline reference. In such case, Eq. 3 simplifies to $(\text{RES} - \text{REL})/\text{UNC}$ (e.g., Mason, 2004), and positive BSS values, therefore, correspond to $\text{RES} > \text{REL}$.

85   Reliability diagrams, also known as Attributes Diagrams (Hsu and Murphy, 1986), are commonly used to illustrate the last two components of the Brier Score. A step-by-step procedure to construct a reliability diagram is as follows. First define

binary event based on a climatological threshold (e.g., value is below the median or above the upper tercile). Second, compute the corresponding climatological thresholds separately for the forecast and verification data at each grid point. Third, define discrete forecast probability bins (for example, 0–0.2, 0.2–0.4, and so on) into which forecast probabilities will be grouped.

90 Fourth, for each year and each grid point within a region (see land regions in Sect. 2.4), compute the forecast probability of the binary event as a fraction of ensemble members for which the binary event occurs, assign this probability to the appropriate bin and record whether the same binary event occurred in the verification data. Finally, within each probability bin, count the number of forecasts that fall into it and the number of times the event occurred in verification data. The observed frequency in probability bin $k$, $\bar{o}_k$, is then calculated as the ratio of the number of occurrences in verification data to the number of forecasts

95 in the bin. Plotting these observed frequencies on vertical axis against the corresponding forecast probabilities on horizontal axis, with the size of each point reflecting the number of forecasts within the corresponding bin, illustrates how well the forecast probabilities match the observed outcomes (Figure 1). In such a diagram, the distance between the $(f_k, \bar{o}_k)$ coordinates and the diagonal represents reliability, while the distance to the horizontal climatological probability line represents resolution.
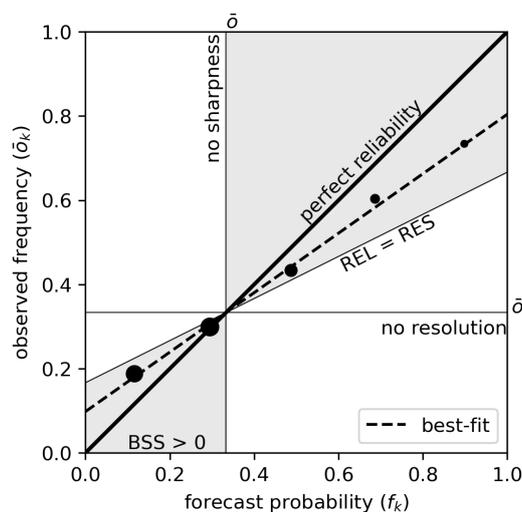


**Figure 1.** An example of a reliability diagram for an upper tercile event (climatological probability $\bar{o} = 1/3$) shown for five forecast probability bins ($k = 1..m$, where $m = 5$). The area with the positive Brier Skill Score (BSS) is shaded gray. The results are based on the ECMWF hindcasts verified against ESA Snow-CCI data over WNAS (Western North Asia) land region (see Sect. 2.4).

Grey area in Figure 1 corresponds to situations where the Brier Skill Score (BSS), i.e. the Brier Score's improvement over

100 no-skill climatology, is positive (see e.g., Hsu and Murphy, 1986; Mason, 2004). Another feature seen in Figure 1 is sharpness which reflects the ability of probabilistic forecasts to predict non-climatological probabilities, that is to issue predictions with probabilities close to 0 and 1. Ideally, a probabilistic forecasting system should provide reliable forecasts spanning a wide range of probability intervals, with as many predictions as possible falling within the lowest (around 0) and highest (around 1) probability bins.

105    Following Weisheimer and Palmer (2014), we fit a weighted linear regression to the data points in the diagram (hereafter best-fit reliability line, shown as dashed line in Figure 1) to estimate an overall reliability. The weights correspond to the number of forecasts within each probability bin, and the slope of the resulting regression line serves as an indicator of reliability. A perfectly reliable forecast would lie along the diagonal line (with the slope of 1), while deviations from this line indicate over-confidence (slope is less than 1) or under-confidence (slope is more than 1). To assess the uncertainty around the best-fit

110    reliability line, we generate 1000 bootstrap resamples with replacement from the forecast–observation pairs along the time dimension, compute slopes of the best-fit reliability lines, and derive the 90 percent confidence interval from the resampling slope distribution. While the results shown in Sect. 3 are based on the i.i.d. bootstrap, their robustness was verified using a block bootstrap approach with the blocks of three years.

       In order to quantify the reliability, we apply categorization based on the slope of the best-fit reliability line and its uncertainty

115    following Weisheimer and Palmer (2014). Five reliability categories can be summarized as following. Category 5 (perfect) - the uncertainty range of the reliability slope lies within the positive BSS area and includes the diagonal, Category 4 (useful) - the uncertainty range of the reliability slope lies within the positive BSS area but does not include the diagonal, Category 3 (marginally useful) - the uncertainty range of the reliability slope is positive but may fall out of the positive BSS area, Category 2 (not useful) - the uncertainty range of the reliability slope includes negative values, Category 1 (dangerous) - the slope of the

120    best-fit reliability line is negative. In addition to these five categories, we further introduce a subdivision for Category 3 based on whether the best-fit reliability slope indicates some degree of skill. Category $3^*$ is hereafter used when the best-fit slope falls within the positive BSS area, whereas Category 3, as defined above, is applied when the best-fit slope does not meet this criterion.

## 2.2    Seasonal snow hindcasts

125    The seasonal hindcasts used in this study are based on an updated configuration of the ECMWF Seasonal forecasting System (SEAS5, Johnson et al., 2019). The major differences to SEAS5 include an updated cycle of the ECMWF Integrated Forecasting System (IFS) with a new multi-layer snow scheme (Arduini et al., 2019). The atmospheric component of the system is the IFS Cycle 48r1 with a spectral resolution of TCO319 (horizontal resolution of ~36 km) and 137 vertical levels extending up to 0.01 hPa. The oceanic component is based on the NEMO (Nucleus for European Modelling of the Ocean) model, version 3.4.

130    Within the CERISE project framework (https://www.cerise-project.eu), four-month-long hindcasts were initialized four times per year (1 February, 1 May, 1 August, 1 November) every year between 1993 and 2022, producing a 30-year dataset. Each hindcast ensemble consists of 51 members. The hindcasts use fifth generation ECMWF reanalysis (ERA5) initial conditions for both the atmosphere and land surface. ERA5 assimilates in-situ measurements of snow depth from the international synoptic network (SYNOP) and space-borne snow cover data from the Interactive Multisensor Snow and Ice Mapping System

135    (IMS) from 2004. IMS snow cover fraction is converted to snow depth using an empirical conversion rule (de Rosnay et al., 2015). Note also that snow DA is restricted to elevations below 1500 m in ERA5 (de Rosnay et al., 2022). The ocean and sea-ice initial conditions are based on the ECMWF ORAS5 reanalysis (Zuo et al., 2019).

In this study, we focus on the winter season (December–January–February, DJF) and analyze a set of SWE hindcasts initialized on 1 November 1993 – 2022.

## 2.3 Verification datasets

In this study, the reliability assessment of snow hindcasts is performed against two SWE products, namely, ERA5-Land reanalysis (Muñoz Sabater et al., 2021) and ESA Snow-CCI version 4 (Luojus et al., 2025).

ERA5-Land is a high-resolution global land-surface reanalysis produced by ECMWF, forced by near-surface atmospheric fields from ERA5. It produces a total of 50 land-surface variables including snow (for the full list of variables, see  Muñoz Sabater et al., 2021). The fine resolution of $0.1°$ makes it particularly valuable for hydrological and climate studies over land. ERA5-Land slightly benefits indirectly from the snow depth and snow temperature DA performed in ERA5 via the atmospheric forcing. The SWE product is provided globally at 1 hour temporal resolution from 1950 to present (Muñoz Sabater, 2019). While hourly data are available, ERA5-Land SWE is extracted at 00 UTC to ensure time consistency with the output from the hindcasts.

The ESA Snow-CCI SWE product is based on a retrieval methodology that combines space-borne observations of passive microwave brightness temperatures and in-situ measurements of snow depth at synoptic weather stations (Pulliainen, 2006; Takala et al., 2011). In the ESA Snow-CCI version 2, used in the SWE benchmarking study by Mudryk et al. (2025), SWE product is post-corrected by accounting for the snow density variability in time and space (Venäläinen et al., 2021). Starting from version 3.1, variability of snow density fields is implemented into the SWE retrieval following Venäläinen et al. (2023). Here, we use the latest ESA Snow-CCI version 4 product. Daily SWE retrievals are provided in the Northern Hemisphere winter period only (between mid-October and mid-May) and are available in years 1979 – 2023. The SWE product is provided at $0.1°$ resolution in the Northern Hemisphere excluding mountainous regions, glaciers and Greenland. The mountain mask is applied to grid cells with high sub-grid elevation variability derived from a high-resolution digital elevation model, where retrievals have known limitations (Luojus et al., 2021; Barella et al., 2024).

## 2.4 Reliability assessment configuration

Seasonal forecasting primarily targets deviations from the climatological mean. Expressing SWE as anomalies automatically removes systematic model biases relative to verification data. The winter-mean SWE anomalies in both ERA5-Land and ESA Snow-CCI products are calculated with respect to the whole 1993 – 2022 period. The hindcast anomalies for each ensemble member are calculated with respect to the model mean over the same period.

We consider two binary events based on terciles of the long-term distribution: i) winter-mean SWE anomaly lies below the lower tercile and ii) winter-mean SWE anomaly lies above the upper tercile at a particular geographical location. In regard to the snow forecasts, these events can be also considered as i) low - and ii) high snow accumulation winters.

Based on the climatology of winter-mean SWE in the hindcasts and verification data, eight land regions spread over the snow covered areas were selected for the reliability assessment (see Table 1 and Fig. 2). The span and names of the land regions are based on land regions used in the study by Giorgi and Francisco (2000) but are adjusted to the needs of our winter

snow analysis. Note that WNA, NEU and TIB regions partially cover areas with low winter-mean SWE ($< 3$ cm) and high year-to-year SWE variability (Figure 2b,d,f).

**Table 1.** List of regions used in this study.

| Name | Acronym | Latitude range | Longitude range | Relative sample size (%) |
|---|---|---|---|---|
| Alaska | ALA | $60°$ N – $72°$ N | $170°$ W – $103°$ W | 47.71 |
| Canada | CND | $45°$ N – $72°$ N | $103°$ W – $60°$ W | 76.15 |
| Western North America | WNA | $35°$ N – $60°$ N | $130°$ W – $103°$ W | 37.29 |
| Norther Europe | NEU | $50°$ N – $75°$ N | $0°$ E – $45°$ E | 61.95 |
| Western North Asia | WNAS | $50°$ N – $75°$ N | $45°$ E – $90°$ E | 99. 89 |
| Central North Asia | CNAS | $50°$ N – $75°$ N | $90°$ E – $135°$ E | 100.00 |
| Eastern North Asia | ENAS | $50°$ N – $75°$ N | $135°$ E – $179°$ E | 41.84 |
| Tibet | TIB | $28°$ N – $50°$ N | $65°$ E – $105°$ E | 53.44 |

To ensure consistency in the reliability assessment, both verification datasets were interpolated onto a uniform latitude–longitude grid with a resolution of $0.25°$ using bilinear interpolation. As mentioned in Sect. 2.3, ESA Snow-CCI SWE product does not cover mountainous regions. To enable a consistent comparison, the ESA Snow-CCI mountain mask, seen as white areas in Figure 2c-d, is applied to ERA5-Land dataset. The reliability assessment is, therefore, performed over the non-mountainous terrain in both cases. In addition, ERA5-Land data points with snow amount exceeding 1 m of SWE (resulting from the interpolation of permanent snow masked as 10 m of SWE) are discarded. This affects regions CND, WNAS and ENAS that include small islands with permanent snow, with the following percentage of data points removed 0.28 %, 0.34 % and 0.22 %, respectively. Note that ESA Snow-CCI has a limit of 0.5 m for SWE product (Luojus et al., 2025). Because water bodies and mountainous terrain are excluded from the analysis, the number of valid land grid points varies across regions. The sample size in each region is expressed as a fraction of the number of grid points in the most populated region (CNAS) as shown in Table 1.

## 3 Results. How reliable are seasonal snow forecasts?

This section summarizes the reliability of ECMWF snow hindcasts in DJF 1993 – 2022 for low- and high snow accumulation winters shown in Figure 3. The overall hindcast performance is good, as only useful categories (Category 3 – 5) are obtained for both verification datasets and tercile events.

Figure 3a shows that assessment against ERA5-Land for low snow accumulation winters yields one region with perfect reliability (Category 5, green), five regions with useful reliability (Category 4, blue), and two regions with marginally useful reliability (Category 3*, yellow). Figure A1 presents individual reliability diagrams for all assessed land regions. The hindcasts show excellent performance in predicting low snow amounts in CND region, although they are not very sharp because relatively few forecasts fall within the high-probability bin. The latter can be explained by the presence of high amount of snow in the
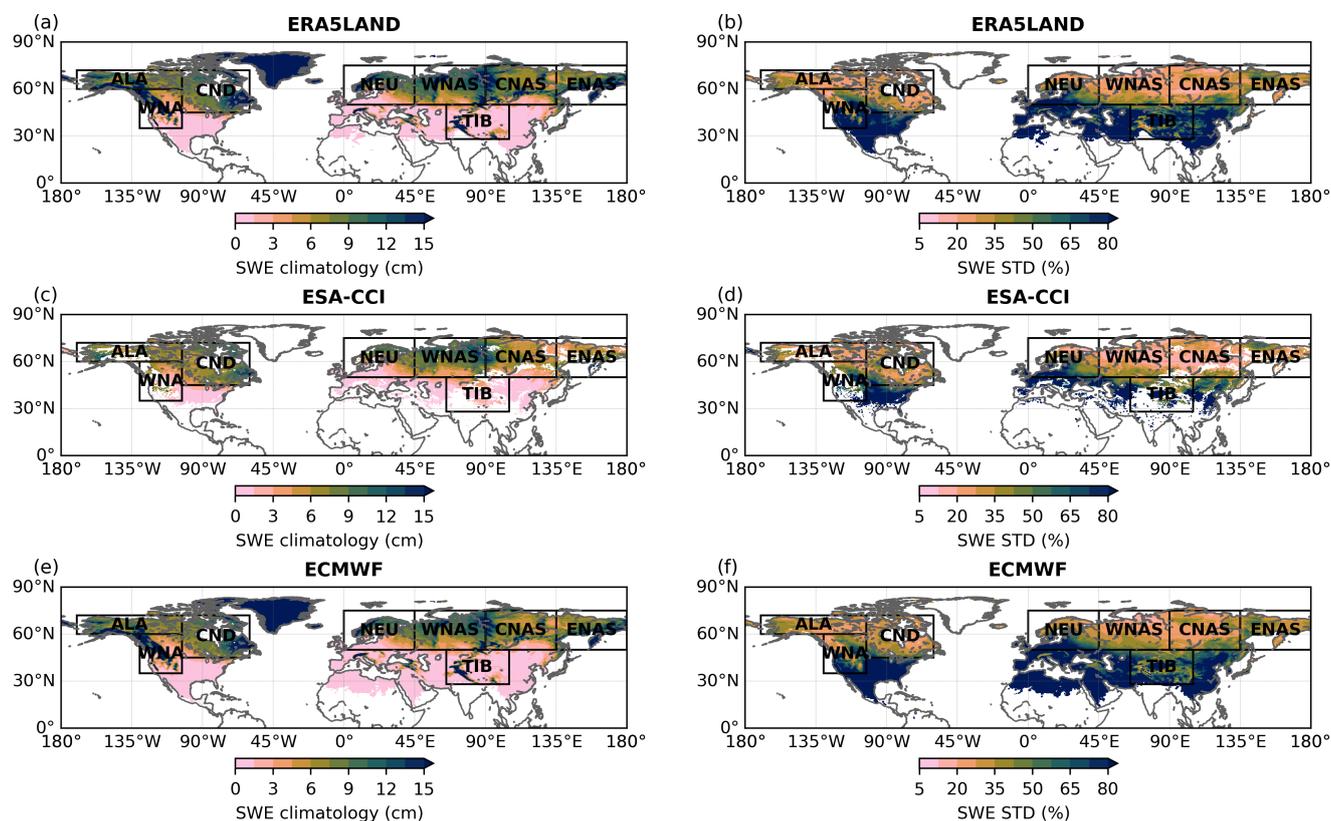
**Figure 2.** Climatological mean SWE and standard deviation (in percent relative to the climatological mean) in DJF 1993 – 2022 in a) – b) ERA5-Land reanalysis, c) – d) ESA Snow-CCI, e) – f) ECMWF hindcasts. Eight land regions in which reliability assessment is performed (namely, ALA, CND, WNA, NEU, WNAS, CNAS, ENAS, TIB) are shown as black boxes.

region most of the time. ALA, WNAS, CNAS, ENAS, and TIB regions also display good skill in predicting low snow amounts. For NEU and WNA regions, which fall into the marginally useful category, most hindcasts issue probabilities close to the climatological value of 1/3, with few forecasts assigning either low or high probabilities. This indicates limited sharpness and resolution in these regions. Nevertheless, their best-fit reliability lines remain within the skillful BSS range and are therefore assigned Category $3^*$.

The results are substantially different when verified against ESA Snow CCI (Figure 3b). In this case, snow hindcasts are assigned marginally useful category in all eight land regions. Although, NEU, WNAS, CNAS and TIB regions show an improved skill since their best-fit reliability lines fall within the skillful BSS range (Category $3^*$, yellow). Figure A2 shows individual reliability diagrams. A clear off-set of the best-fit reliability line from the climatological intersection ($f_k = \bar{o}$, $\bar{o}_k = \bar{o}$) is present in TIB and WNA regions. Such situation reflects unconditional bias in the hindcasts compared to the verification data (Weisheimer and Palmer, 2014). It was mentioned in Sect. 2.4 that these regions partially cover snow transition areas with

high year-to-year SWE variability. The biases between ERA5, ERA5-Land, in-situ and satellite observations over the Tibetan

205  Plateau (encompassed by the region labeled TIB here) have been documented by Orsolini et al. (2019).
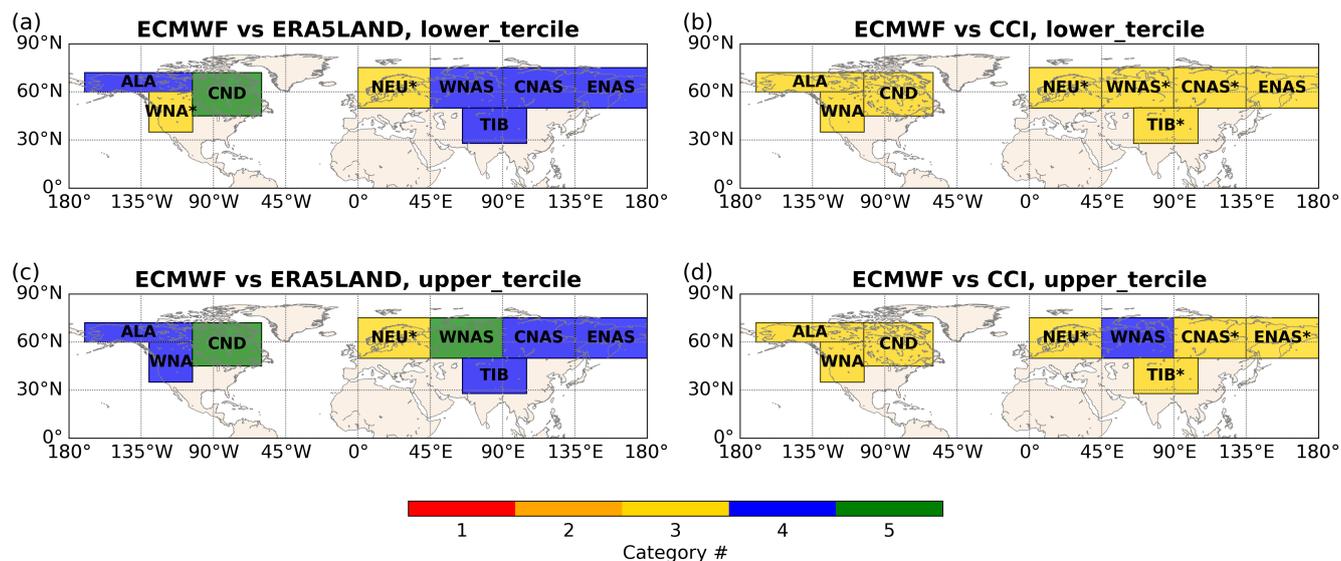


**Figure 3.** Reliability of the ECMWF snow hindcasts in DJF 1993 – 2022 for low snow accumulation (lower tercile) winters assessed against a) ERA5-Land, b) ESA Snow-CCI, and for high snow accumulation (upper tercile) winters assessed against c) ERA5-Land, d) ESA Snow-CCI. Reliability categories are color coded. Note that a star sign next to the region name indicates Category 3$^*$ (i.e. when the best-fit slope falls within the BSS > 0 area) as described in Sect. 2.4.

When verified against ERA5-Land in high snow accumulation winters, snow hindcasts are assigned perfect (Category 5, green) reliability category in two land regions and useful (Category 4, blue) category in five land regions (Figure 3c), and marginally useful (Category 3$^*$, yellow) category in NEU region. Individual reliability diagrams are shown in Figure A3. Similarly to the low snow accumulation winters, the hindcasts lack sharpness and resolution in NEU region. In addition, the

210  best-fit reliability line is on the edge of falling inside the skillful BSS area, which suggests that forecasts lack skill and do not perform much better than climatology (see Sect. 2.1 for explanation).

When verifying against ESA Snow-CCI, the overall picture is different – only one out of eight land regions is assigned useful (Category 4, blue) reliability category while other seven regions are assigned marginally useful category (Category 3, yellow), see Figure 3d. Although, NEU, CNAS, ENAS and TIB regions show an improved skill and are therefore assigned Category

215  3$^*$. Individual reliability diagrams are shown in Figure A4. Similar to the low snow accumulation winters, TIB region shows a clear bias between the SWE in hindcasts and ESA Snow-CCI data.

## 4 Discussion

Figure 3 demonstrates that reliability assessment and categorization are somewhat dependent on the choice of verification dataset. Although it may seem intuitive to place greater trust in Earth observation products, these datasets have notable limitations. It is known that ESA Snow-CCI has a systematic bias in estimates of large SWE values ($> 15$ cm) due to the reduced sensitivity of passive-microwave retrievals to deep snow and liquid water in the snowpack (Luojus et al., 2021; Venäläinen et al., 2025). In addition, the mismatch between point-scale in-situ measurements being assimilated and the SWE grid-scale estimates adds uncertainty into the final product, with retrieval performance declining farther from assimilated snow depth observations (Luojus et al., 2021). A recent study by Venäläinen et al. (2025) suggests bias correction for ESA Snow-CCI SWE retrievals that could be used in the future assessments. At the same time, ERA5-Land, while generally robust (Mudryk et al., 2025), also has known biases and difficulties in SWE estimation (Monteiro and Morin, 2023; Kouki et al., 2023).

Based on the individual reliability diagrams shown in Figures A1 – A4, the ECMWF hindcasts tend to be over-confident i.e., they issue probabilities that are too extreme (too low for low probability bins and too high for high probability bins). Data assimilation plays an important role in constraining the forecast initial snow to observations. The type of assimilated data and the assimilation method are, therefore, of a key importance. In the hindcasts used here, both land and atmosphere initial conditions are derived from ERA5 reanalysis. Since the off-line ERA5-Land reanalysis uses ERA5 as atmospheric forcing, higher reliability categories obtained when verifying against ERA5-Land may partly reflect the shared dependence on ERA5. Consistent with this interpretation, similarly high categories were found when the hindcasts were verified against ERA5 itself (not shown). The results should, therefore, be interpreted as conditional on this shared ERA5 forcing dependence, rather than as a standalone measure of forecast reliability.

To better understand differences in reliability categories between the two verification datasets, Figure 4 shows probability density distributions of SWE anomalies in ERA5-Land and ESA Snow-CCI in different land regions. Despite similarities in probability density distributions, it is clear that SWE anomalies in ERA5-Land have more extreme values and exhibit more variability, while anomalies in ESA Snow-CCI are more closely clustered around the mean. Time series of snow anomalies in Figure 5 also show discrepancies in their temporal evolution and in the distribution of years with strong positive and negative anomalies. The latter is important for the definition of probabilistic tercile thresholds as explained in Sect. 2.1. Figure A5 compares distributions of lower- and upper tercile thresholds in ERA5-Land and ESA Snow-CCI. The two datasets generally show similar distributions of terciles, however with shifts in the peaks of their distributions. Differences in distributions strongly vary between regions. TIB region (non-mountainous) stands out with much narrower tercile distributions in ERA5-Land centered close to zero compared to smoother tercile distributions in ESA Snow-CCI.
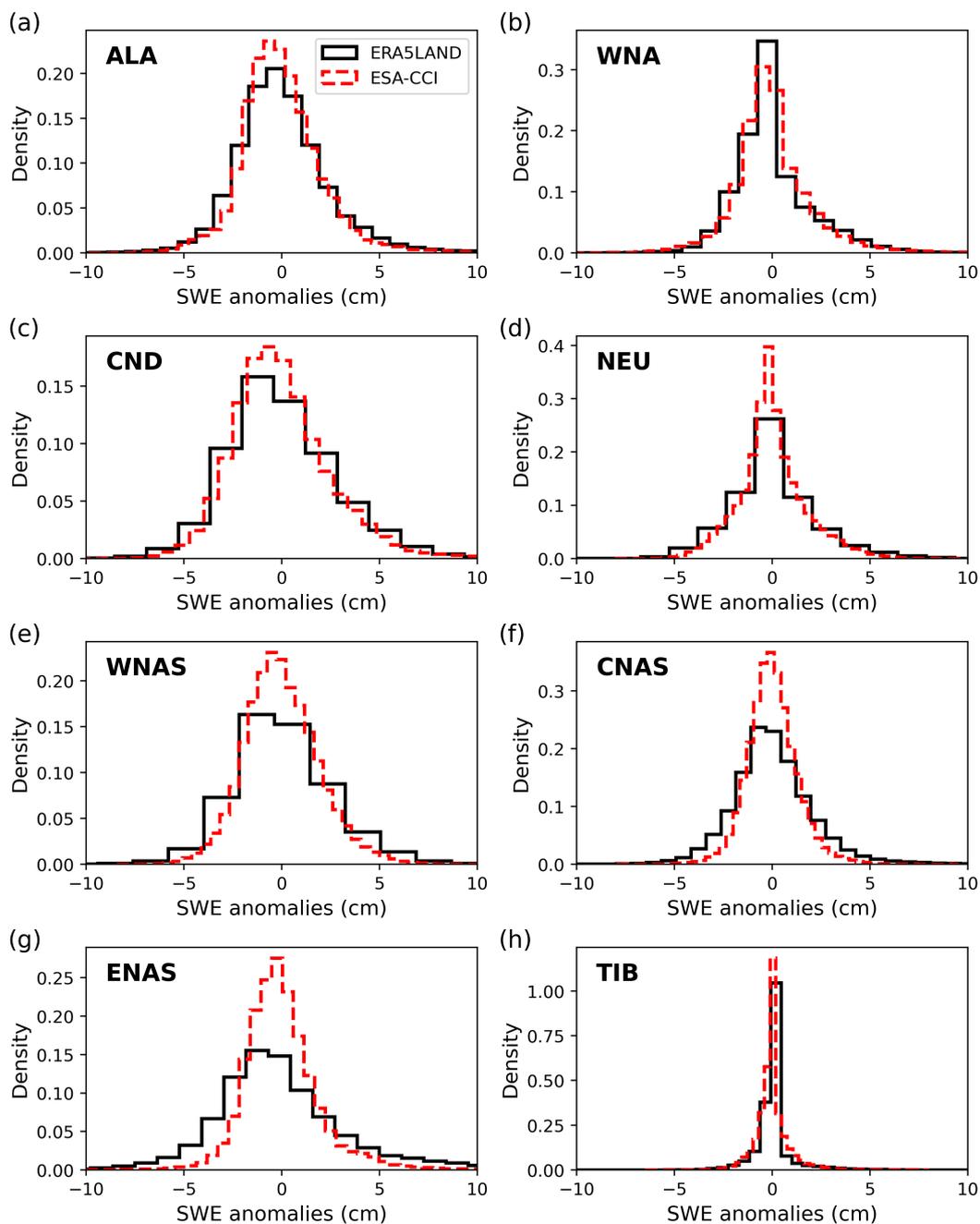
**Figure 4.** Probability density distributions of SWE anomalies in ERA5-Land (black solid) and ESA Snow-CCI (red dashed) in eight land regions: a) ALA, b) WNA, c) CND, d) NEU, e) WNAS, f) CNAS, g) ENAS, h) TIB. Note that scale in panels is different.
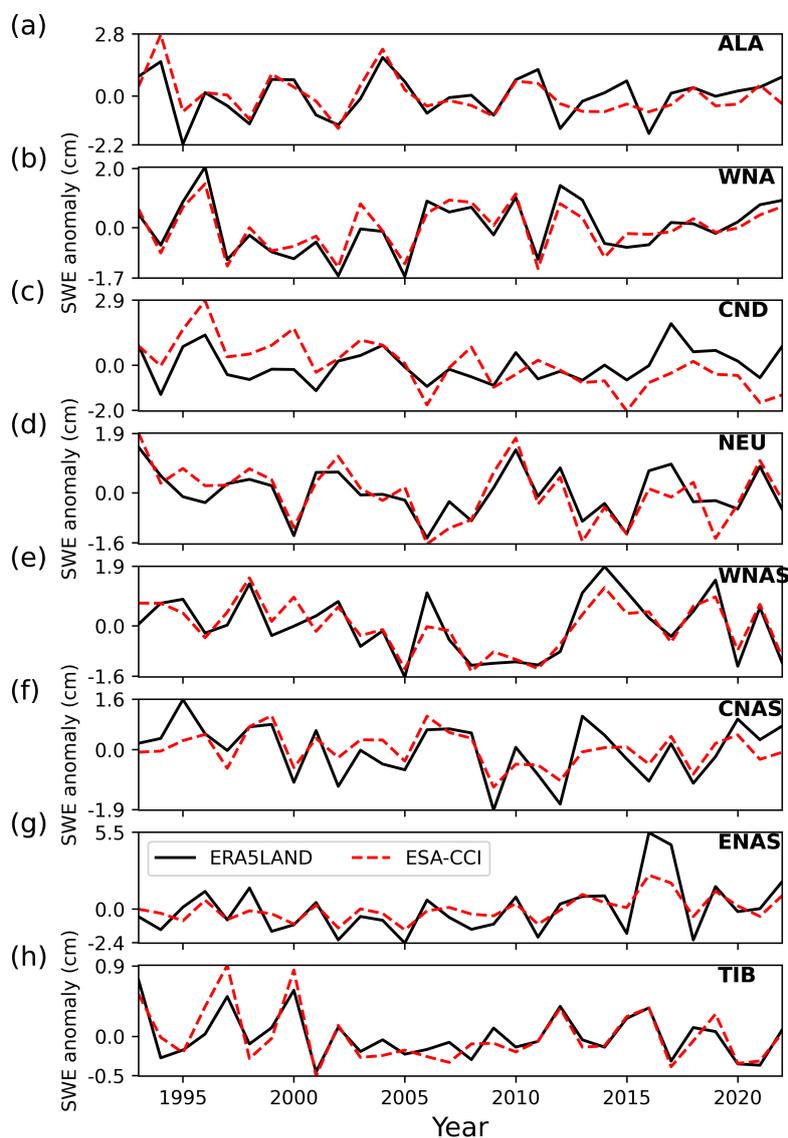
**Figure 5.** Time series of SWE anomalies in ERA5-Land (black solid) and ESA Snow-CCI (red dashed) averaged in land regions: a) ALA, b) WNA, c) CND, d) NEU, e) WNAS, f) CNAS, g) ENAS, h) TIB. Note that scale in panels is different.

It is also important to examine how often the two validation datasets agree on the binary event classification, i.e., whether it occurs or not. Table A1 shows the percentage of data points where the binary event is classified consistently in both datasets at the same geographical location and time. Overall agreement is computed using all data points within the region. Near-threshold agreement is computed using only those data points for which SWE anomalies in both validation datasets lie within a distance of 1 cm of SWE from their respective tercile thresholds at the same location and time. While the overall agreement between

250

**12**

ERA5-Land and ESA Snow-CCI in terms of binary event classification varies around 70–80%, the agreement for cases near the tercile threshold drops to approximately 50–70%. This indicates substantial differences in event classification close to the threshold, which directly affect the categorical reliability assessment and explain the discrepancies observed in the reliability diagrams shown in Figures A1 – A4.

255   As noted in Sect. 1, numerous gridded snow products are currently available, however those frequently show substantial discrepancies in both magnitude and spatial–temporal variability of snow, as well as in long-term snow trends (see e.g., Mudryk et al., 2025). Such inconsistencies emphasize the critical need for continued efforts to develop a high-quality, internally consistent snow dataset specifically suited for the forecast verification purposes.

## 5   Conclusions

260   In this study, we assess reliability of winter-mean snow hindcasts in 1993 – 2022 produced by the ECMWF within the CERISE project. In probabilistic forecasting, reliability for a binary event is defined as the consistency between forecast probabilities and observed frequencies. The assessment is performed against two SWE datasets: ERA5-Land reanalysis and ESA Snow-CCI (version 4) in eight non-mountainous land regions in the Northern Hemisphere. To evaluate the forecast performance in low and high snow accumulation winters, we consider two binary events based on terciles of the long-term distribution: i) winter-mean

265   SWE anomaly lies below the lower tercile, and ii) winter-mean SWE anomaly lies above the upper tercile.

A simple categorization is used to quantify the reliability based on the weighted linear regression to the data in reliability diagrams. The ECMWF snow hindcasts show an overall good performance as only high reliability categories (Category 3 – 5) are obtained for both verification datasets and tercile events. Although the reliability assessment is sensitive to the choice of verification dataset and yields slightly different reliability categories when verified against ERA5-Land and ESA Snow-

270   CCI, this sensitivity itself highlights the current limitations of two best available SWE datasets (Mudryk et al., 2025). Such inconsistencies emphasize the critical need for continued efforts to develop a high-quality, internally consistent snow dataset specifically suited for the forecast verification purposes.

Answering the question "How reliable are seasonal forecasts of snow?", we show that the ECMWF hindcasts exhibit good reliability in predicting winter-mean snow water equivalent in both low- and high-accumulation winters. The reliability vari-

275   ations among the regions could be linked to, for example, regional differences in snow distribution: high SWE variability in snow transition areas increases the challenge of capturing the full range of possible outcomes, which can indirectly lead to lower forecast reliability. In addition, technical limitations in both the forecast model (e.g., dependence on the empirical snow cover-snow depth conversion rule) and the verification datasets (e.g., the aforementioned ESA Snow-CCI bias for large SWE values) can have an effect on the spatial distribution of the hindcasts reliability.

280   Avenues for future work include, but are not limited to, a more detailed assessment of seasonal SWE forecast reliability as a function of lead time, for example, through monthly assessments throughout the winter season. In addition, exploring prediction of extreme snow events, like the snow droughts mentioned in Sect. 1, would be of a great interest as they have been linked to hydrological extremes. It is also important to investigate the reliability of SWE forecasts in mountainous areas in

spring, when snowmelt occurs with potential implications for hydropower management and flood-risk assessment. Finally, a

285  multi-model comparison of seasonal SWE forecasts produced by operational meteorological prediction centers would provide

additional insight into model performance and reliability.

*Data availability.* Hourly ERA5 data on single levels from 1940 to present are available from the Climate Data Store website (Hersbach et al., 2023) via https://doi.org/10.24381/cds.adbb2d47. ESA-CCI SWE v4 data are available from the CEDA website (Luojus et al., 2025) via https://dx.doi.org/10.5285/edf8abd23f4a40aabd4d52e48dec06ea. ECMWF hindcasts are available via the Meteorological Archival and

290  Retrieval System (MARS) under "CERISE project" dataset.

*Author contributions.* Ekaterina Vorobeva and Yvan Orsolini initiated the study and wrote manuscript with contributions from Patricia de Rosnay, Jonathan Day, Retish Senan, Frederic Vitart, and Damien Decremer. Ekaterina Vorobeva performed calculations and made figures.

*Competing interests.* Authors declare that no competing interests are present.
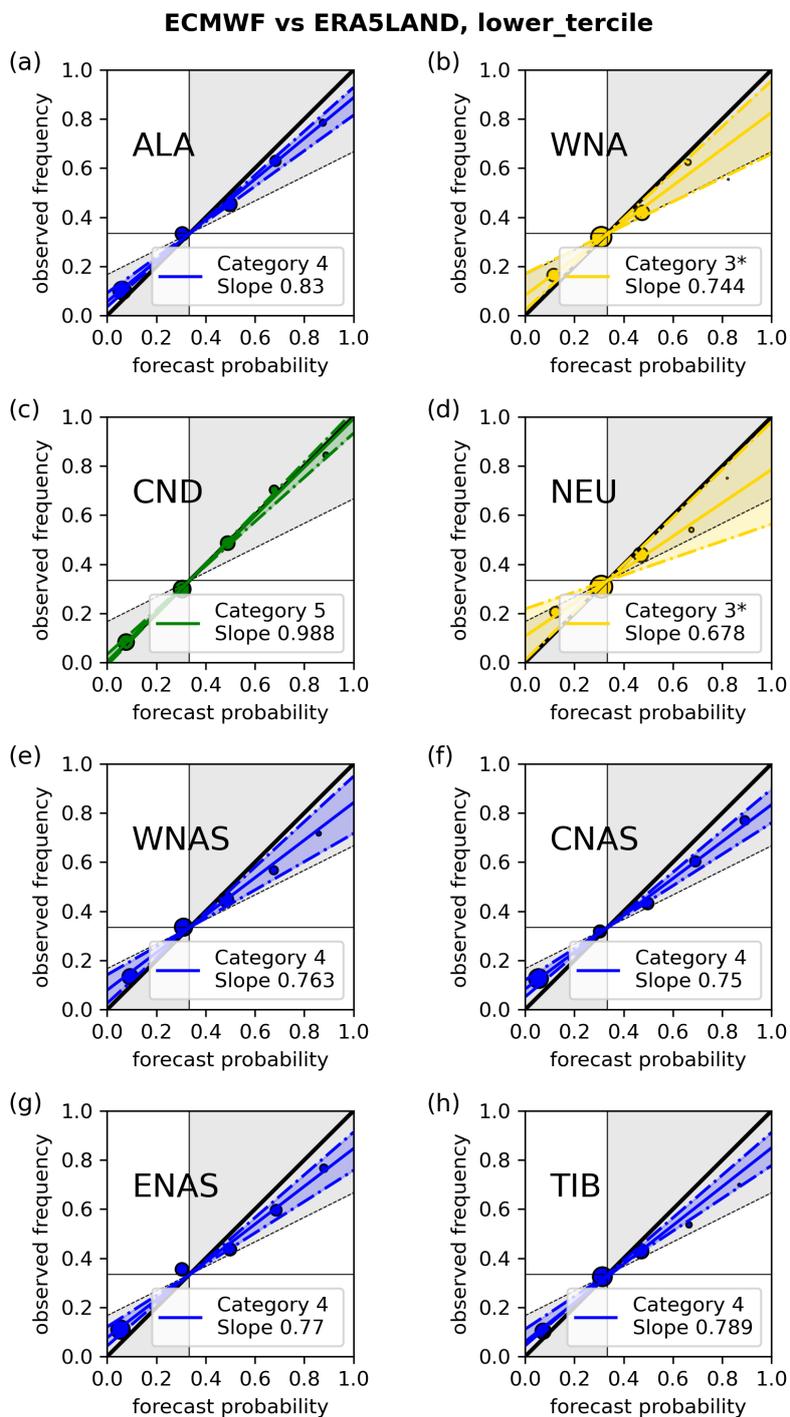
**Appendix A**

**Figure A1.** Reliability diagrams for the ECMWF SWE hindcasts validated against ERA5-Land in low snow accumulation winters over eight land regions: a) ALA, b) WNA, c) CND, d) NEU, e) WNAS, f) CNAS, g) ENAS, h) TIB.
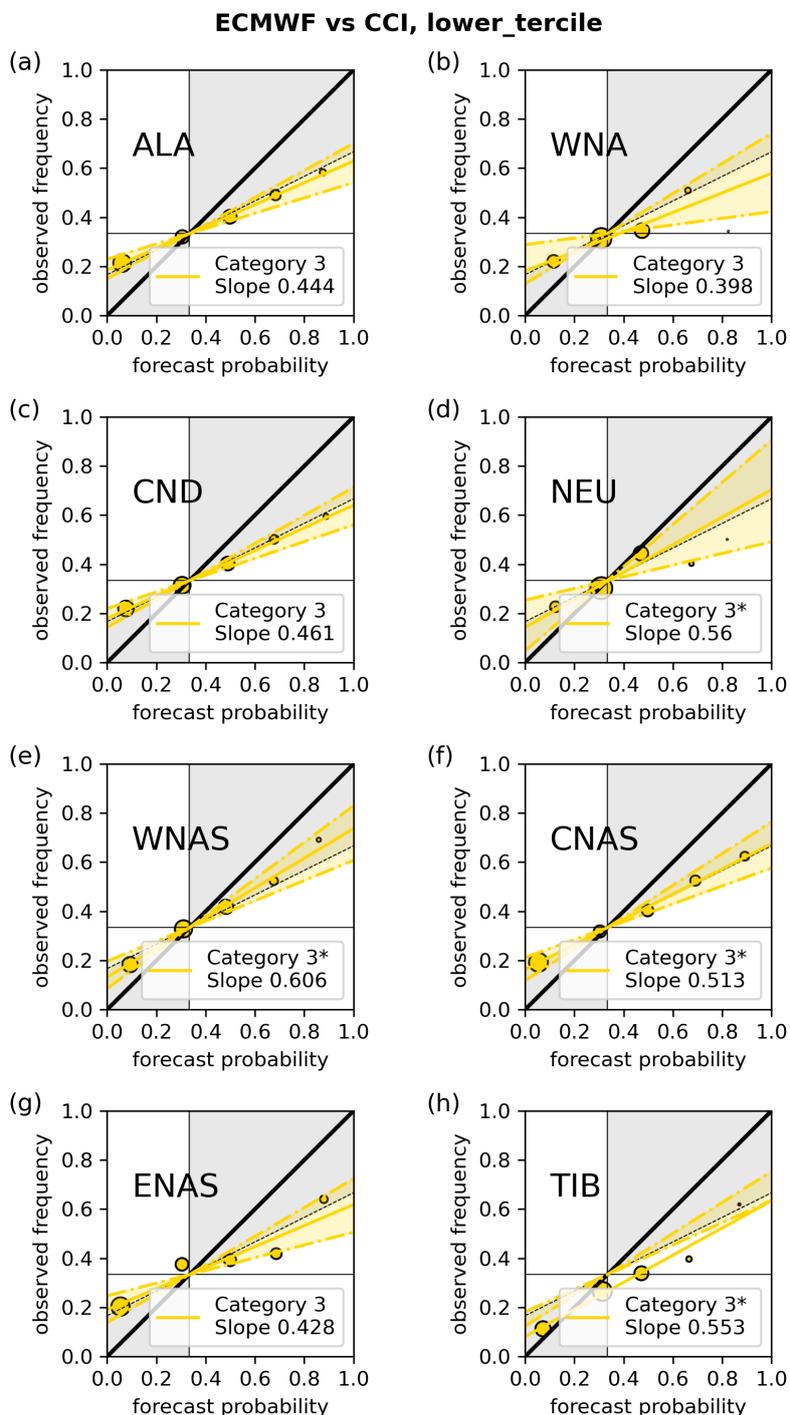
**ECMWF vs CCI, lower_tercile**



**Figure A2.** Reliability diagrams for the ECMWF SWE hindcasts validated against ESA Snow-CCI in low snow accumulation winters over eight land regions: a) ALA, b) WNA, c) CND, d) NEU, e) WNAS, f) CNAS, g) ENAS, h) TIB.
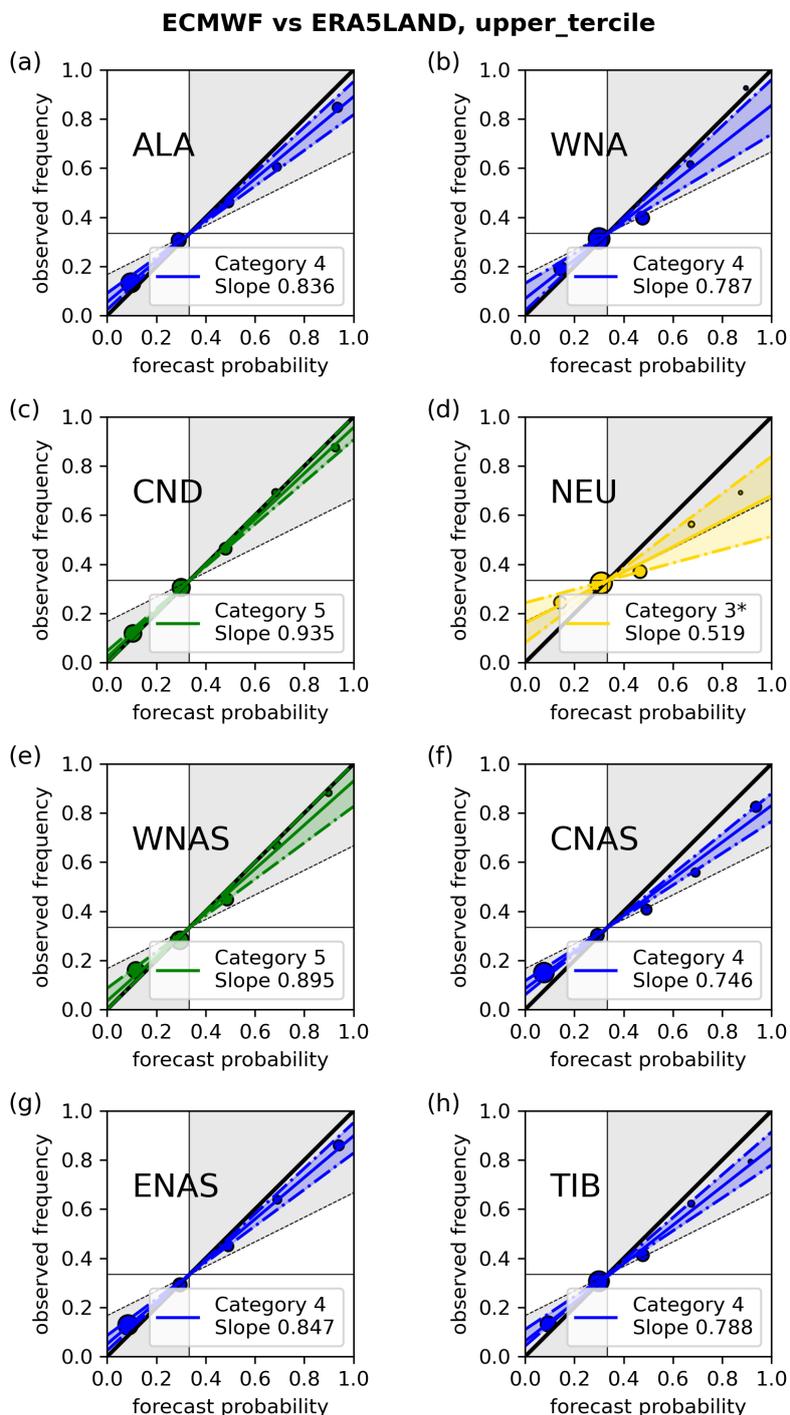
**Figure A3.** Reliability diagrams for the ECMWF SWE hindcasts validated against ERA5-Land in high snow accumulation winters over eight land regions: a) ALA, b) WNA, c) CND, d) NEU, e) WNAS, f) CNAS, g) ENAS, h) TIB.
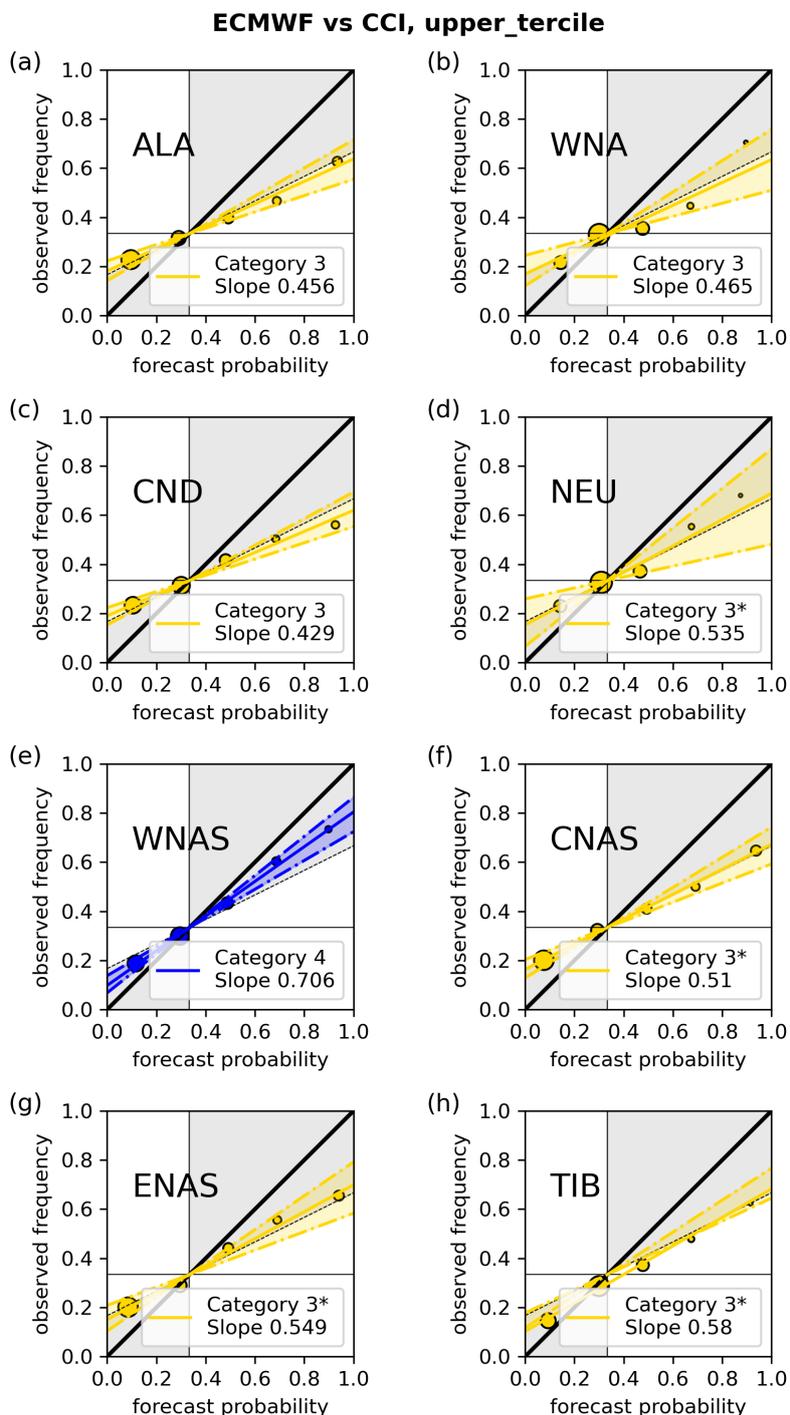
**Figure A4.** Reliability diagrams for the ECMWF SWE hindcasts validated against ESA Snow-CCI in high snow accumulation winters over eight land regions: a) ALA, b) WNA, c) CND, d) NEU, e) WNAS, f) CNAS, g) ENAS, h) TIB.
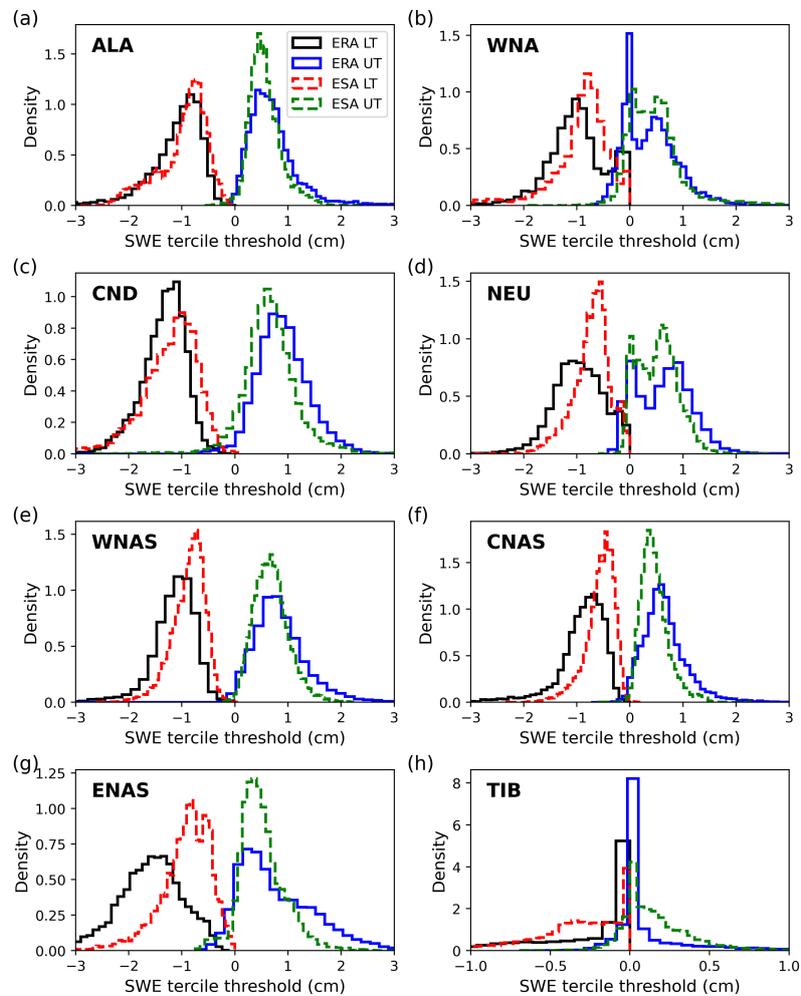
**Figure A5.** Probability density distributions of lower- and upper terciles (LT and UT, respectively) in ERA5-Land (black solid for LT, blue solid for UT) and ESA Snow-CCI (red dashed for LT, green dashed for UT) in eight land regions: a) ALA, b) WNA, c) CND, d) NEU, e) WNAS, f) CNAS, g) ENAS, h) TIB.

| Region | Agreement in % (overall / near-threshold) |
|:------:|:----------------------------------------:|
| ALA | 71 / 54 |
| WNA | 80 / 70 |
| CND | 69 / 54 |
| NEU | 80 / 71 |
| WNAS | 80 / 59 |
| CNAS | 76 / 61 |
| ENAS | 74 / 54 |
| TIB | 74 / 71 |

**Table A1.** The summary of agreement between ERA5-Land and ESA Snow-CCI in classifying binary events, i.e. percentage of data points where binary event occurs in both datasets. Overall agreement is evaluated over all data points within the region. Agreement near the tercile threshold is restricted to cases where SWE anomalies satisfy $|anomaly - T| < w$ in both datasets, T denotes the tercile-based threshold of the corresponding dataset and w is selected as 1 cm. The condition is applied point-wise to ensure that both datasets have near-threshold SWE values at the same space–time location.

# References

300    Arduini, G., Balsamo, G., Dutra, E., Day, J. J., Sandu, I., Boussetta, S., and Haiden, T.: Impact of a Multi-Layer Snow Scheme on Near-Surface Weather Forecasts, Journal of Advances in Modeling Earth Systems, 11, 4687–4710, https://doi.org/https://doi.org/10.1029/2019MS001725, 2019.

Barella, R., Mortimer, C., Marin, C., Schwaizer, G., Mölg, N., Nagler, T., Wunderle, S., Xiao, X., Luojus, K., Venäläinen, P., Takala, M., Pulliainen, J., Lemmetyinen, J., Moisander, M., and Solberg, R.: ESA CCI+ Snow ECV: Product User Guide, version 4.0, 2024.

305    Bormann, K. J., Brown, R. D., Derksen, C., and Painter, T. H.: Estimating snow-cover trends from space, Nature Climate Change, 8, 924–928, https://doi.org/10.1038/s41558-018-0318-3, 2018.

Brier, G. W.: Verification of forecasts expressed in terms of probability, Monthly weather review, 78, 1–3, 1950.

de Rosnay, P., Isaksen, L., and Dahoui, M.: Snow data assimilation at ECMWF, https://doi.org/10.21957/lkpxq6x5, 2015.

de Rosnay, P., Browne, P., de Boisséson, E., Fairbairn, D., Hirahara, Y., Ochi, K., Schepers, D., Weston, P., Zuo, H., Alonso-Balmaseda,
310    M., et al.: Coupled data assimilation at ECMWF: current status, challenges and future developments, Quarterly Journal of the Royal Meteorological Society, 148, 2672–2702, https://doi.org/10.1002/qj.4330, 2022.

Fox-Kemper, B., Hewitt, H., Xiao, C., Aðalgeirsdóttir, G., Drijfhout, S., Edwards, T., Golledge, N., Hemer, M., Kopp, R., Krinner, G., et al.: Ocean, cryosphere and sea level change. Climate change 2021: the physical science basis. Contribution of Working Group I to the Sixth Assessment Report of the Intergovernmental Panel on Climate Change, Climate Change, pp. 1211–1362,
315    https://doi.org/10.1017/9781009157896.011, 2021.

Giorgi, F. and Francisco, R.: Uncertainties in regional climate change prediction: a regional analysis of ensemble simulations with the HADCM2 coupled AOGCM, Climate Dynamics, 16, 169–182, https://doi.org/10.1007/PL00013733, 2000.

Hersbach, H., Bell, B., Berrisford, P., Biavati, G., Horányi, A., Muñoz Sabater, J., Nicolas, J., Peubey, C., Radu, R., Rozum, I., Schepers, D., Simmons, A., Soci, C., Dee, D., and Thépaut, J.-N.: ERA5 hourly data on single levels from 1940 to present [data set]. Copernicus
320    Climate Change Service (C3S) Climate Data Store (CDS), https://doi.org/10.24381/cds.adbb2d47, accessed on 28-OCT-2025, 2023.

Hsu, W.-r. and Murphy, A. H.: The attributes diagram a geometrical framework for assessing the quality of probability forecasts, International Journal of Forecasting, 2, 285–293, https://doi.org/10.1016/0169-2070(86)90048-8, 1986.

Jeong, J.-H., Linderholm, H. W., Woo, S.-H., Folland, C., Kim, B.-M., Kim, S.-J., and Chen, D.: Impacts of snow initialization on subseasonal forecasts of surface air temperature for the cold season, Journal of Climate, 26, 1956–1972, https://doi.org/10.1175/JCLI-D-12-00159.1,
325    2013.

Johnson, S. J., Stockdale, T. N., Ferranti, L., Balmaseda, M. A., Molteni, F., Magnusson, L., Tietsche, S., Decremer, D., Weisheimer, A., Balsamo, G., Keeley, S. P. E., Mogensen, K., Zuo, H., and Monge-Sanz, B. M.: SEAS5: the new ECMWF seasonal forecast system, Geoscientific Model Development, 12, 1087–1117, https://doi.org/10.5194/gmd-12-1087-2019, 2019.

Komatsu, K. K., Takaya, Y., Toyoda, T., and Hasumi, H.: A submonthly scale causal relation between snow cover and surface air temperature
330    over the autumnal Eurasian continent, Journal of Climate, 36, 4863–4877, https://doi.org/10.1175/JCLI-D-22-0827.1, 2023.

Kouki, K., Luojus, K., and Riihelä, A.: Evaluation of snow cover properties in ERA5 and ERA5-Land with several satellite-based datasets in the Northern Hemisphere in spring 1982–2018, The Cryosphere, 17, 5007–5026, https://doi.org/10.5194/tc-17-5007-2023, 2023.

Li, F., Orsolini, Y., Keenlyside, N., Shen, M.-L., Counillon, F., and Wang, Y.: Impact of snow initialization in subseasonal-to-seasonal winter forecasts with the Norwegian Climate Prediction Model, Journal of Geophysical Research: Atmospheres, 124, 10033–10048,
335    https://doi.org/10.1029/2019JD030903, 2019.

Luojus, K., Pulliainen, J., Takala, M., Lemmetyinen, J., Mortimer, C., Derksen, C., Mudryk, L., Moisander, M., Hiltunen, M., Smolander, T., et al.: GlobSnow v3. 0 Northern Hemisphere snow water equivalent dataset, Scientific Data, 8, 163, https://doi.org/10.1038/s41597-021-00939-2, 2021.

Luojus, K., Venäläinen, P., Moisander, M., Pulliainen, J., Takala, M., Lemmetyinen, J., Mortimer, C., Mudryk, L., Schwaizer, G., and Nagler,
340    T.: ESA Snow Climate Change Initiative (Snow_cci): Snow Water Equivalent (SWE) level 3C daily global climate research data package (CRDP) (1979 - 2023), version 4.0 [data set], https://doi.org/https://dx.doi.org/10.5285/edf8abd23f4a40aabd4d52e48dec06ea, accessed on 28-OCT-2025, 2025.

Mason, S. J.: On using "climatology" as a reference strategy in the Brier and ranked probability skill scores, Monthly Weather Review, 132, 1891–1895, https://doi.org/10.1175/1520-0493(2004)132%3C1891:OUCAAR%3E2.0.CO;2, 2004.

345    Monteiro, D. and Morin, S.: Multi-decadal analysis of past winter temperature, precipitation and snow cover data in the European Alps from reanalyses, climate models and observational datasets, The Cryosphere, 17, 3617–3660, https://doi.org/10.5194/tc-17-3617-2023, 2023.

Muñoz Sabater, J., Dutra, E., Agustí-Panareda, A., Albergel, C., Arduini, G., Balsamo, G., Boussetta, S., Choulga, M., Harrigan, S., Hersbach, H., Martens, B., Miralles, D. G., Piles, M., Rodríguez-Fernández, N. J., Zsoter, E., Buontempo, C., and Thépaut, J.-N.: ERA5-Land: a state-of-the-art global reanalysis dataset for land applications, Earth System Science Data, 13, 4349–4383, https://doi.org/10.5194/essd-
350    13-4349-2021, 2021.

Mudryk, L., Santolaria-Otín, M., Krinner, G., Ménégoz, M., Derksen, C., Brutel-Vuilmet, C., Brady, M., and Essery, R.: Historical Northern Hemisphere snow cover trends and projected changes in the CMIP6 multi-model ensemble, The Cryosphere, 14, 2495–2514, https://doi.org/10.5194/tc-14-2495-2020, 2020.

Mudryk, L., Mortimer, C., Derksen, C., Elias Chereque, A., and Kushner, P.: Benchmarking of snow water equivalent (SWE) products based
355    on outcomes of the SnowPEx+ Intercomparison Project, The Cryosphere, 19, 201–218, https://doi.org/10.5194/tc-19-201-2025, 2025.

Muñoz Sabater, J.: ERA5-Land hourly data from 1950 to present [data set]. Copernicus Climate Change Service (C3S) Climate Data Store (CDS), https://doi.org/10.24381/cds.e2161bac, accessed on 28-OCT-2025, 2019.

Notarnicola, C.: Hotspots of snow cover changes in global mountain regions over 2000–2018, Remote Sensing of Environment, 243, 111 781, https://doi.org/10.1016/j.rse.2020.111781, 2020.

360    Ombadi, M., Risser, M. D., Rhoades, A. M., and Varadharajan, C.: A warming-induced reduction in snow fraction amplifies rainfall extremes, Nature, 619, 305–310, https://doi.org/10.1038/s41586-023-06092-7, 2023.

Orsolini, Y., Senan, R., Balsamo, G., Doblas-Reyes, F., Vitart, F., Weisheimer, A., Carrasco, A., and Benestad, R.: Impact of snow initialization on sub-seasonal forecasts, Climate dynamics, 41, 1969–1982, https://doi.org/10.1007/s00382-013-1782-0, 2013.

Orsolini, Y., Wegmann, M., Dutra, E., Liu, B., Balsamo, G., Yang, K., de Rosnay, P., Zhu, C., Wang, W., Senan, R., et al.: Evaluation of
365    snow depth and snow cover over the Tibetan Plateau in global reanalyses using in situ and satellite remote sensing observations, The Cryosphere, 13, 2221–2239, https://doi.org/10.5194/tc-13-2221-2019, 2019.

Pulliainen, J.: Mapping of snow water equivalent and snow depth in boreal and sub-arctic zones by assimilating space-borne microwave radiometer data and ground-based observations, Remote sensing of Environment, 101, 257–269, https://doi.org/10.1016/j.rse.2006.01.002, 2006.

370    Takala, M., Luojus, K., Pulliainen, J., Derksen, C., Lemmetyinen, J., Kärnä, J.-P., Koskinen, J., and Bojkov, B.: Estimating northern hemisphere snow water equivalent for climate research through assimilation of space-borne radiometer data and ground-based measurements, Remote Sensing of Environment, 115, 3517–3529, https://doi.org/10.1016/j.rse.2011.08.014, 2011.

Venäläinen, P., Luojus, K., Lemmetyinen, J., Pulliainen, J., Moisander, M., and Takala, M.: Impact of dynamic snow density on GlobSnow snow water equivalent retrieval accuracy, The Cryosphere, 15, 2969–2981, https://doi.org/10.5194/tc-15-2969-2021, 2021.

375 Venäläinen, P., Luojus, K., Mortimer, C., Lemmetyinen, J., Pulliainen, J., Takala, M., Moisander, M., and Zschenderlein, L.: Implementing spatially and temporally varying snow densities into the GlobSnow snow water equivalent retrieval, The Cryosphere, 17, 719–736, https://doi.org/10.5194/tc-17-719-2023, 2023.

Venäläinen, P., Mortimer, C., Luojus, K., Mudryk, L., Takala, M., and Pulliainen, J.: Updated monthly and new daily bias correction for assimilation-based passive microwave SWE retrieval, The Cryosphere, 19, 6301–6318, https://doi.org/10.5194/tc-19-6301-2025, 2025.

380 Weisheimer, A. and Palmer, T. N.: On the reliability of seasonal climate forecasts, Journal of The Royal Society Interface, 11, 20131 162, https://doi.org/10.1098/rsif.2013.1162, 2014.

Zhang, W., Liu, L., Wu, H., Zhang, T., Chen, Y., and Wang, L.: Snow droughts amplify compound climate extremes over the Tibetan Plateau, Communications Earth & Environment, 6, 571, https://doi.org/10.1038/s43247-025-02551-3, 2025.

Zuo, H., Balmaseda, M. A., Tietsche, S., Mogensen, K., and Mayer, M.: The ECMWF operational ensemble reanalysis–analysis system
385 for ocean and sea ice: a description of the system and assessment, Ocean Science, 15, 779–808, https://doi.org/10.5194/os-15-779-2019, 2019.