# THE RISE OF DATA-DRIVEN WEATHER FORECASTING
## A FIRST STATISTICAL ASSESSMENT OF MACHINE LEARNING-BASED WEATHER FORECASTS IN AN OPERATIONAL-LIKE CONTEXT

### A PREPRINT

Zied Ben Bouallègue, Mariana C A Clare, Linus Magnusson, Estibaliz Gascón, Michael Maier-Gerber, Martin Janoušek, Mark Rodwell, Florian Pinault, Jesper S Dramsch, Simon T K Lang, Baudouin Raoult, Florence Rabier, Matthieu Chevallier, Irina Sandu, Peter Dueben, Matthew Chantry, Florian Pappenberger

ECMWF

### ABSTRACT

Data-driven modeling based on machine learning (ML) is showing enormous potential for weather forecasting. Rapid progress has been made with impressive results for some applications. The uptake of ML methods could be a game-changer for the incremental progress in traditional numerical weather prediction (NWP) known as the "quiet revolution" of weather forecasting. The computational cost of running a forecast with standard NWP systems greatly hinders the improvements that can be made from increasing model resolution and ensemble sizes. An emerging new generation of ML models, developed using high-quality reanalysis datasets like ERA5 for training, allow forecasts that require much lower computational costs and that are highly-competitive in terms of accuracy. Here, we compare for the first time ML-generated forecasts with standard NWP-based forecasts in an operational-like context, initialized from the same initial conditions. Focusing on deterministic forecasts, we apply common forecast verification tools to assess to what extent a data-driven forecast produced with one of the recently developed ML models (PanguWeather) matches the quality and attributes of a forecast from one of the leading global NWP systems (the ECMWF IFS). The results are very promising, with comparable skill for both global metrics and extreme events, when verified against both the operational analysis and synoptic observations. Increasing forecast smoothness and bias drift with forecast lead time are identified as current drawbacks of ML-based forecasts. A new NWP paradigm is emerging relying on inference from ML models and state-of-the-art analysis and reanalysis datasets for forecast initialization and model training.

## 1 Introduction

Numerical weather prediction (NWP) is the dominant approach for weather forecasting. A weather forecast is the result of the numerical integration of partial differential equations starting from the best estimate of the current state of the Earth System. The idea that the physical laws of fluid dynamics and thermodynamics can be used to predict the state of the atmosphere dates back to the pioneering works of Abbe (1901) and Bjerknes (1904). In a standard NWP framework, a weather prediction results from a deductive inference: a deterministic forecast is derived using the laws of physics starting from the best possible initial conditions, derived by optimally combining earth system observations and short-range forecasts through data assimilation. However, our ability to perfectly know the initial conditions and numerically resolve the equations is limited. Hence, ensemble forecasting is used to account for uncertainty in both the initial conditions and the forecasting model, with the resulting ensemble forecast serving as a basis for probabilistic forecasting (Leutbecher and Palmer, 2008).

A continuous improvement of the NWP performance has been observed over the last decades, including for the prediction of high-impact weather events (Ben Bouallègue et al., 2019). Skill improvement is achieved through improvements in initial conditions, numerical models, and resolution. At the European Centre for Medium-range Weather Forecasting (ECMWF), the Integrated Forecasting System (IFS) has been run operationally since 1979 with regular updates of the different components of the forecasting system. The evolution of the IFS skill over the last two

decades is shown in Fig. 1 (red lines). The steady increase in forecast skill thanks to incremental improvements in numerical modeling, supercomputing, data assimilation and ensemble techniques, observations and their use in the NWP system, has become known as the 'quiet revolution' of weather forecasting (Bauer et al., 2015). However, the computational cost of running a forecast is a major bottleneck that hinders rapid improvements with standard NWP systems. In operational NWP, the computational and timeliness constraints imply finding a balance between increasing model resolution and increasing ensemble size, which are two major factors known to improve the skill of ensemble forecasts (Leutbecher and Ben Bouallègue, 2020).

Recently, data-driven modeling based on machine learning (ML) is showing large potential for weather forecasting applications (de Burgh-Day and Leeuwenburg, 2023). Tremendous progress has been made since 2022 with a series of key works developing machine-learning models for weather forecasting (Keisler, 2022; Pathak et al., 2022; Bi et al., 2023; Lam et al., 2022; Chen et al., 2023b). These recent publications present impressive forecast scores, some of which rival the operational ECMWF deterministic high-resolution (deterministic) forecasts. Concretely, Keisler (2022) use a Graph Neural Network (GNN) model and claim to produce more accurate forecasts of specific humidity than IFS after day 3; Pathak et al. (2022) leverage Fourier Transforms with a transformer and claim to produce comparable accuracy to IFS for 2m temperature; Bi et al. (2023) use a vision transformer model and claim to produce more accurate forecasts than IFS across numerous variables when both models are verified against reanalysis; Lam et al. (2022) use a GNN and claim more accurate forecasts than IFS on a larger set of atmospheric variables and pressure levels; and finally Chen et al. (2023b) use a transformer and claim to improve the scores compared to Lam et al. (2022), especially at longer lead times.

The emergence of data-driven models has been made possible thanks to the availability of large high-quality meteorological free and open datasets. The aforementioned ML models are trained on ERA5 reanalysis data, which is the fifth generation ECMWF atmospheric reanalysis, produced by the Copernicus Climate Change service, as one of the European Union Copernicus Programme key deliverables (Hersbach et al., 2020). This dataset is particularly attractive for machine learning problems because it is a continuous weather dataset from 1940 to the present day and it represents the best possible reconstruction of the Earth-system state created by blending past observations and short-range forecasts through data assimilation. However, the ML methods presented only train from 1979 because the extensions to 1940 are relatively recent and have lower accuracy due to the very limited availability of satellite data before 1980 (e.g. Hersbach, 2023). ERA5 is generated using the operational IFS cycle at the time of production (2016) and is publicly available at a grid resolution of $0.25°$ (30km). Hence, ML models are trained on reanalysis at a much lower resolution than that of today's operational forecasts and analyses (30km instead of 9km in the case of the ECMWF operational high-resolution forecasts and analyses). Note that, despite this resolution difference, ERA5 hindcasts[1] are used routinely for forecast verification purposes: the performance of the current IFS is compared with the performance of ERA5 forecasts (10-day forecast initialized from ERA5 at the ERA5 resolution, about 25 km) to help distinguish inter-annual variability from actual skill improvement due to changes in the forecasting system. This is illustrated in Fig. 1 with ERA5 hindcast skill represented by black lines.

One approach to data-driven weather prediction would consist of running ML-trained models starting from optimized initial conditions in an operational context. In such a weather prediction system, the forecast inference relies on an ML model rather than the physical model (included for example in the IFS). This approach is highly attractive because a forecast can be generated at a speed several orders of magnitude faster than that from conventional methods. At a fundamental level, an ML-based prediction is the result of an inductive rather than a deductive inference. This paradigm shift in terms of logic has implications for the way a weather forecast is interpreted: a forecast becomes a plausible outcome given what has been learned from previous data. However, the mode of inference followed by ML methods can raise concerns, in particular regarding the ability of such models to predict extreme events unseen in the training dataset. Moreover, the interpretability of ML models is also often questioned when they are perceived as black boxes where the link between the training dataset and the current forecast is difficult to grasp (McGovern et al., 2019). The huge potential benefits and drawbacks of data-driven systems trigger the question of whether ML models can become a component of operational NWP systems.

In this study, we evaluate the performance of data-driven forecasts in an operational-like context. More precisely, the PanguWeather ML model in Bi et al. (2023) (referred to hereafter as PGW), which is open-source for non-commercial use, has been set up to run on the ECMWF computers. For the first time, a forecast generated with an ML model is compared with an operational NWP forecast using the same framework and starting from the same initial conditions. In Bi et al. (2023), like in previous studies, the ML-based forecasts were initialized from ERA5. We leverage standard verification techniques routinely applied for weather forecast evaluation at ECMWF. Using this methodology, we can assess which aspects of the data-driven forecasts can match the quality of forecasts performed with one of the leading operational NWP systems. This study focuses predominantly on the statistical analysis of forecast performance, but we

---

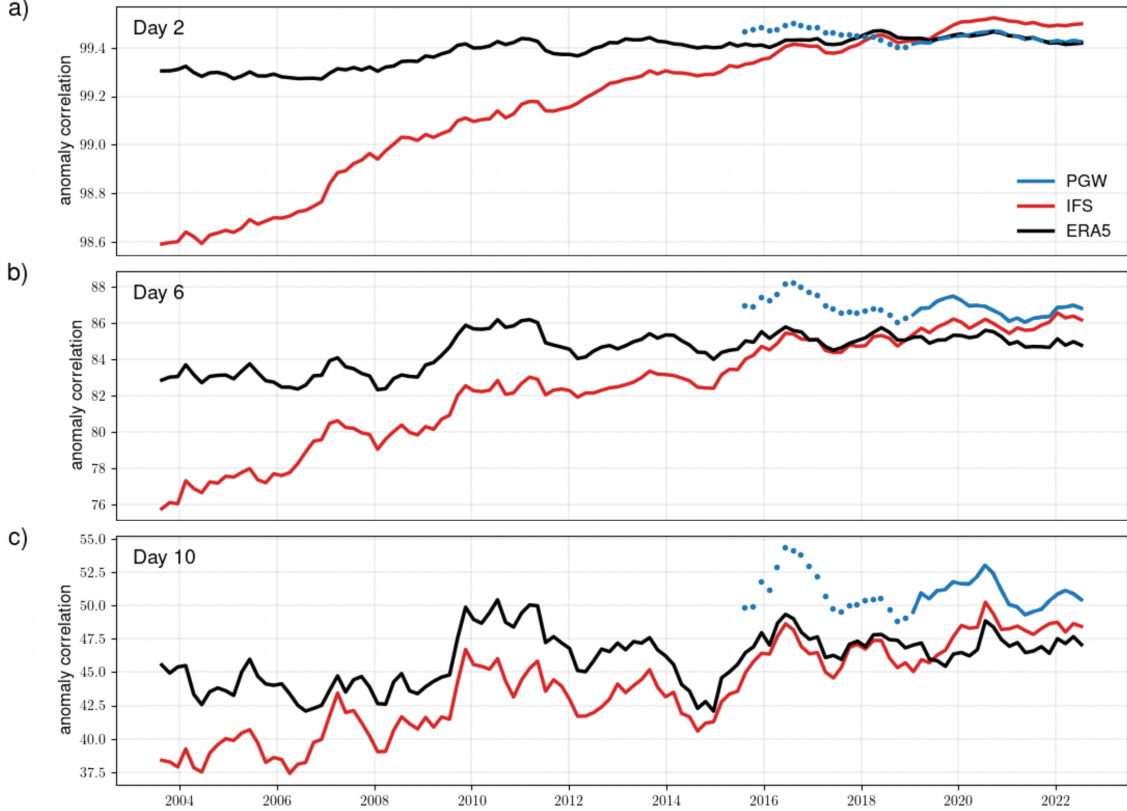[1]forecasts for a historical period

Figure 1: Forecast skill (the larger the better) over the Northern Hemisphere at day 2 (a), day 6 (b), and day 10 (c). Skill is measured as the correlation between the forecasts and the verifying analysis for the geopotential height at 500-hPa, expressed as the anomaly with respect to the climatological height. A one-year running mean is applied. The constant improvement over the past decades of the IFS forecasts is compared with the performance of the ERA5 hindcasts (run with the IFS version operational in 2016) and with the performance of the PGW hindcasts trained over 1979-2018 and verified over 2015-2023 (the period of overlap between training/validation and verification is shown with dotted lines).

acknowledge that case studies play a key role in understanding the capability and limitations of ML models in weather forecasting and refer the reader to Magnusson (2023).

## 2   Methodology and experiments

Our comparative work is based on implementing PGW in an operational-like setting. PGW uses a vision transformer model architecture (Dosovitskiy et al., 2020) with 3D weather fields as inputs and outputs. Developed in Bi et al. (2023), the network minimizes a loss function defined as the root mean square error (RMSE) with a cos-latitude weighting to account for the spherical nature of the Earth as the model is trained on a regular latitude-longitude grid. Like most ML models, PGW uses an iterative method to forecast forward in time. A novelty of their approach, however, is that they choose to minimize the RMSE over a series of fixed short time periods (1hr, 3hrs, 6hrs, and 24hrs) and then achieve weather forecasts at any time using a Hierarchical Temporal Aggregation method. Here, PGW is run for 10 days, using a 24h-step. Though PGW demonstrated stability beyond this time period, only the first 10 days are analyzed here.

For the verification periods under focus in this work, IFS is run operationally at a horizontal grid resolution of 9km up to 10 days lead time using the operational IFS cycle at the time (43r3 and 45r1 for 2018, 47r3 for 2022). We also include the publicly available ERA5 hindcasts in our comparison. ERA5 hindcasts start from ERA5 reanalysis and are based on a lower model resolution than the operational IFS forecast (30km instead of 9km). Here we recall that the ERA5 reanalysis and hindcasts are produced with a similar set-up of IFS as that used for the operational high-resolution forecasts and analysis, but they are produced with the operational cycle at the time when the production of ERA5 started

(41r2, in 2016), and at a lower resolution (30km instead of 9km). Thus, ERA5 hindcasts and IFS forecasts differ in terms of initial conditions, resolution, and IFS cycle.

We choose to initialize PGW from the same analysis as IFS, namely the ECMWF operational analysis. This choice appears natural for a fair comparison between PGW and IFS in an operational-like setting. PGW was trained using ERA5 data, meaning initialization from the operational analysis may induce some impact on scores. An optimal configuration of an ML-based forecasting system would likely 'fine-tune' (training near convergence) on operational analysis. This optimization is outside the scope of this work but is worthy of note. The operational analysis is created by using the current operational data assimilation system, which operates at a higher resolution (9km) and uses a more recent (and therefore improved) IFS version than ERA5 to construct superior initial conditions.

As a complementary experiment, we also run PGW starting from ERA5 analysis (PGW_E5). This experiment shows that PGW starting from the operational high-resolution analysis generally performs better than PGW starting from ERA5 for the first days of the forecast (roughly up to day 4 depending on the variable and domain of interest), as illustrated in Fig. 2.

We also run IFS initialized from the operational analysis at a lower resolution (IFS_LR) close to the PGW resolution to isolate the impact of model horizontal resolution on forecast performance. This impact differs depending on both variable and lead time. For T850 at a lead time of two days, a change in horizontal resolution results in a clear degradation of the forecast skill, but at a 6 day lead time and for Z500 there are only small differences between IFS and IFS_LR errors in our results in Fig. 2(c).

As expected, IFS_LR is ranked between ERA5 and IFS in terms of performance. Indeed, IFS and IFS_LR forecasts are better than ERA5 forecasts because they start from the operational analysis. Concerning PGW, which has the same horizontal resolution as IFS_LR, it is interesting to note that both forecasts have similar errors to ERA5 for T850 over the winter period, and these are noticeably larger than with IFS operational forecasts. Also, the impact of the model horizontal resolution is smaller at longer lead times, as illustrated in Figs 2(c,d).

## 3 Data and a case-study

We assess the performance of two upper-air variables, geopotential at 500 hPa (Z500) and temperature at 850 hPa (T850). Z500 and T850 forecasts are verified against the IFS operational analysis interpolated to a grid resolution of 1.5° following the World Meteorological Organisation (WMO) guideline and aggregated over the Northern Hemisphere. We also assess forecasts of 2m temperature against surface synoptic observations (SYNOP) over Europe. In addition, a verification of tropical cyclone (TC) forecasts is performed. Details about the verification process are provided in the appendix.

We show results mainly for two seasons: Summer 2022 (1 June 2002 to 31 August 2022) and Winter 2022/2023 (1 December 2022 to 28 February 2023) to allow a focus on both extreme cold and extreme warm temperatures. These two verification periods are independent of the PWG training/validation dataset. Only results for Winter 2022 are shown for the upper variables because it is the most dynamically active season in the northern hemisphere. Following Bi et al. (2023), the TC verification period covers 2 January to 30 November 2018.

A comparison against SYNOP observations helps demonstrate the forecast performance from a user perspective. Nevertheless, in-situ observations have their drawbacks: the quality of the measurements is not perfect, the stations are not distributed homogeneously over the verification domain, and measurements can suffer discontinuity at a given station. Also, representativeness is a major concern when comparing model output with a point observation that might not be representative of the surrounding area. This representativeness issue is partially addressed here by an orography correction applied to the 2m temperature forecasts (see for example Ingleby, 2015).

An illustration of a forecast and observation is provided in Fig. 3 where we include the ensemble forecast run operationally at ECMWF (ENS). The 50-member ensemble has a horizontal-grid resolution of 18km for the period considered here. The forecast evolution over consecutive starting times shows that the ensemble spread becomes smaller as we approach the observation date. In Sondankylä (Finland), -29°C was observed on that occasion. The PGW forecast has an earlier hint of the event severity than the IFS forecast, but both overestimated the temperature significantly to a similar degree. In the corresponding maps, PGW forecasts appear smoother than IFS forecasts, deprived of smaller scale structures, both on day 2 and day 6. This first subjective assessment of a single case study agrees with the statistical analysis of the forecast performance discussed in the next section.
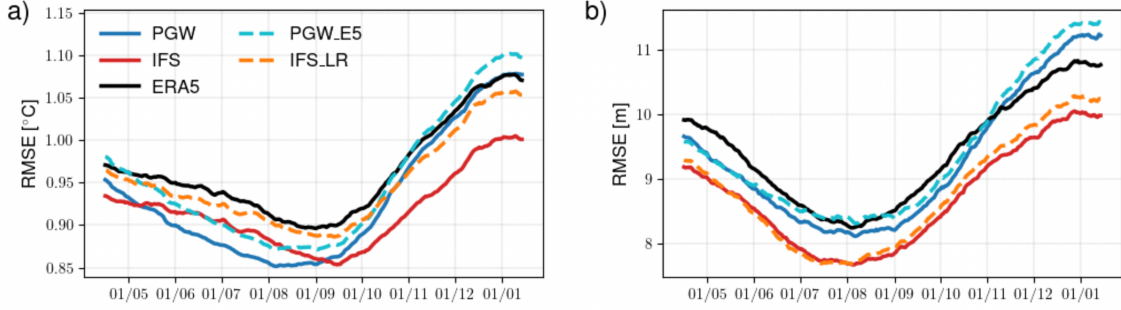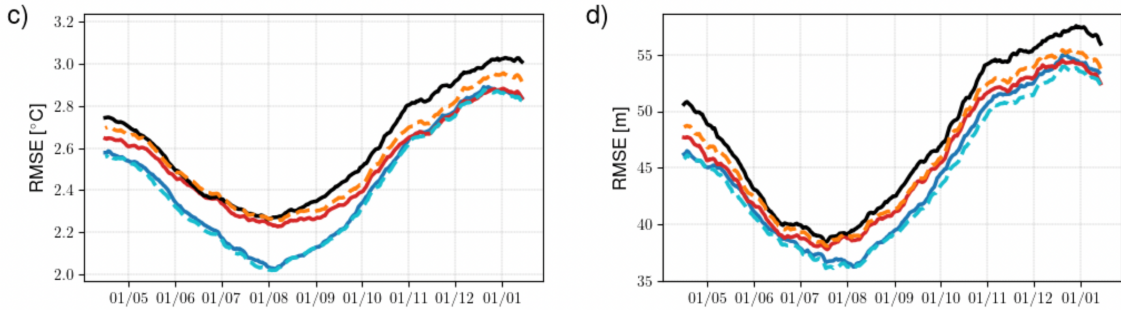
**Day 2**



**Day 6**



Figure 2: The seasonal cycle of the RMSE (3-month averaged) at day 2 and day 6 over one year covering the period 01/03/2022 to 28/02/2023 aggregated over the Northern Hemisphere for (a,c) T850 and (b,d) Z500. PGW_E5: PGW initialized with ERA5, IFS_LR: IFS run at a grid-resolution of 0.25°. The forecasts are verified against operational analyses on a grid resolution of 1.5°.

## 4 Comparing forecast performance

### 4.1 Contextualising the forecast skill

Results in Fig. 4 (top row) are compelling: for lead times greater than 3 days, PGW forecasts are better than the ERA5 forecasts and as good as the operational IFS forecasts in terms of RMSE. The ensemble mean, the ensemble functional that minimises the RMSE, is also a better forecast according to this metric. RMSE is a key indicator of forecast performance but RMSE results need to be interpreted in light of other forecast characteristics that also contribute to the quality of a prediction, for example the realism of the atmospheric state.

In previous studies, there has been a concern that training towards RMSE results in overly smooth forecast fields (see for example the smooth forecasts shown in Keisler, 2022). Indeed, the RMSE strongly penalizes large forecast departures from the observations (or analyses), thus discouraging bold forecasts. When comparing RMSE from different models, it is therefore important to check the level of activity of the different forecasts while interpreting the results. IFS and ERA5 forecasts have similar activity to each other and, importantly, similar activity to PGW for both T850 and Z500 (Fig. 4, middle row).

However, a clear smoothing of PGW forecasts at small scales is visible in Fig. 3. This dampening only has a minor contribution to the overall activity because this metric is dominated by larger scales. We note that for Z500, the activity of PGW slightly decreases with lead time but in general PGW does not become smoother at longer lead times as confirmed by power spectra analysis (not shown). This is not the case for the ensemble mean (EM), which becomes smoother as the forecast uncertainty increases. The lower activity of EM contributes to this good RMSE performance at longer time ranges, but forecast smoothness is a non-desirable characteristic for some applications.

Finally, Fig. 4 also compares the bias of the different forecasts. Ideally, a forecast should have a bias close to zero. The magnitude of the bias in PGW forecasts grows at a much faster rate than the bias in IFS or ERA5 forecasts, with the bias
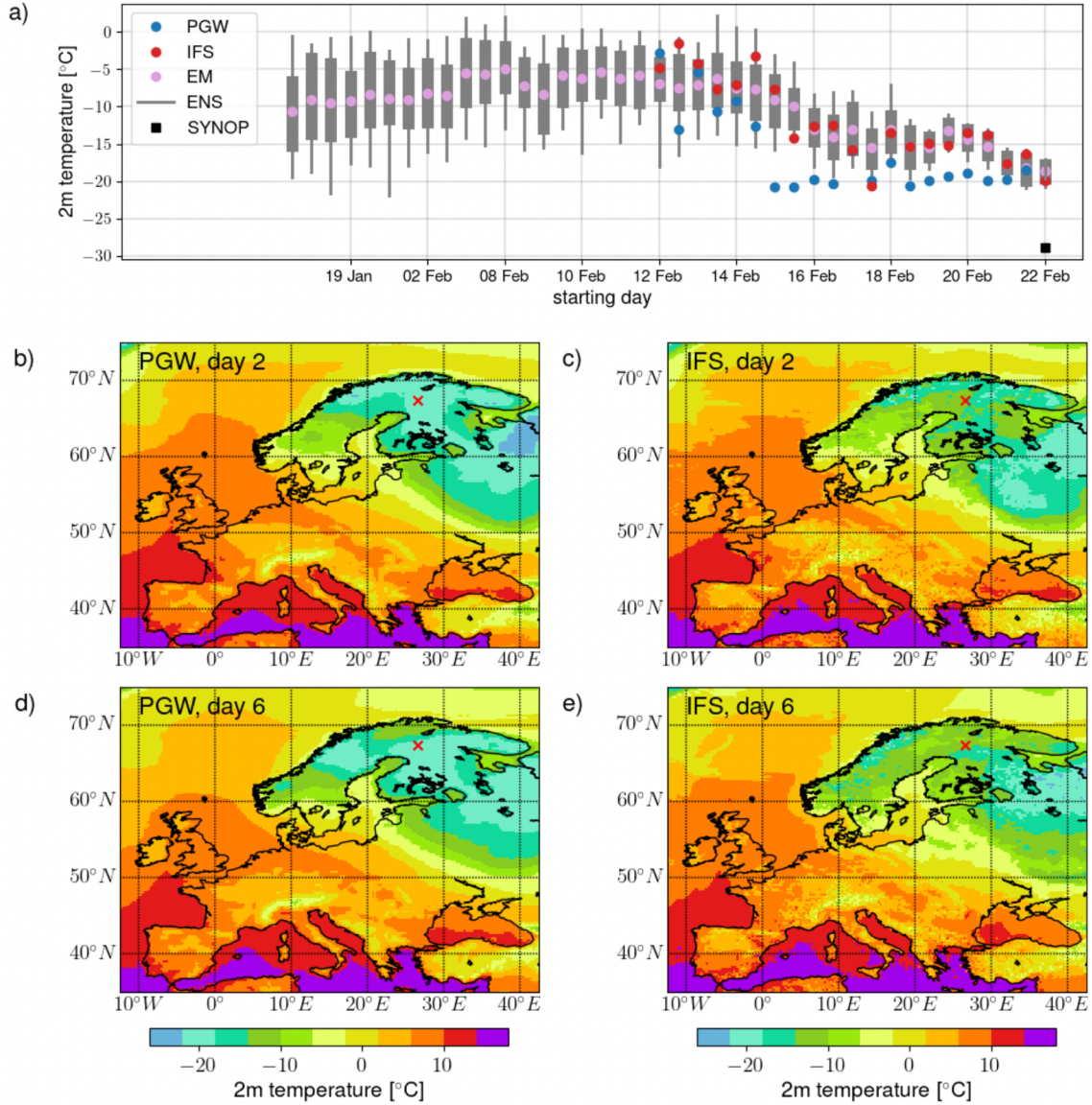
Figure 3: An example of 2m temperature forecasts and a corresponding SYNOP observation. (a) Evolution plots showing forecasts valid at Sodankylä (Finland) on 22 February 2022: the 15-day ensemble forecast in the form of the ensemble mean and quantile forecasts (box-plots showing the 5%, 25%, 75%, and 95%), the 10-day PGW and IFS forecasts. (b) PGW forecast at day 2, (c) IFS forecast at day 2, (d) PGW forecast at day 6, (e) IFS forecast at day 6 for Europe with the location of the Sodankylä SYNOP station indicated with a red cross.

drift particularly strong for Z500. Whereas the bias of IFS and ERA5 stabilizes at longer lead times, the incremental bias in PGW forecasts is still present when extending the forecast horizon.

Verification of upper variables against analysis is complemented by verification of 2m temperature against SYNOP observations in Fig. 5. We find that verifying against observations shows similar results as against analysis for key metrics: the good performance of PGW forecasts in terms of RMSE, a bias drift with forecast lead time in summer, and a larger bias than the IFS and ERA forecasts up to day 5 in winter. Now, verification against observations is used as a framework for a more in-depth analysis of PGW forecast attributes.
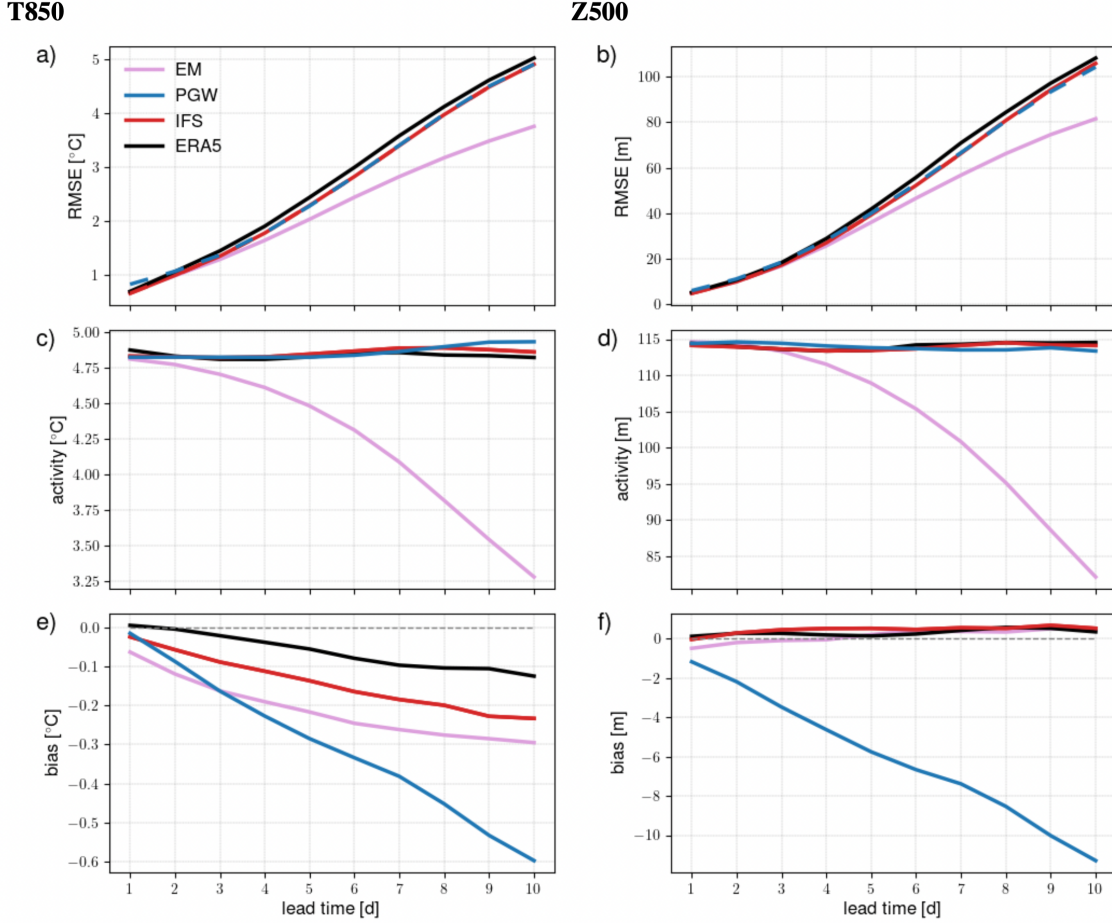
**T850**  **Z500**



Figure 4: (a,b) RMSE (the lower the better), (c,d) forecast activity, (the lower the forecast activity the smoother the forecast), and (e,f) forecast bias (the closer to zero the better), as a function of the forecast lead time for T850 (left panels) and Z500 (right panels). The forecasts are verified against the operational analyses, and results are valid for Winter 2022/2023 over the Northern Hemisphere.

## 4.2  Checking for statistical consistency

Statistical consistency is a key attribute of a forecast. Here we try to answer questions like "Is the forecast able to mimic the observation statistical distribution?", "Is the forecast able to forecast extreme events of the same intensity as the observed ones?" and "Is the forecast systematically offset with respect to the observations?". Statistical consistency in terms of distribution is analysed using quantile-quantile (Q-Q) plots and observation rank histograms. The coherence of the spatial structures in the forecast would require additional diagnostic tools beyond the scope of this study.

Q-Q plots for forecasts at day 6 focus on warm temperatures during summer in Fig. 6(a) and on cold temperatures during winter in Fig. 6(c). For the former, both PGW and IFS forecasts can capture the observed extreme temperatures with PGW displaying a general offset consistent with the PGW bias at day 6 in Fig. 5. For the latter, extremely low temperatures are not fully captured, neither by PGW nor by IFS, as already illustrated in the case study in Fig. 3. In Northern Europe, very low temperatures are reached closest to the ground during clear-sky nights over snow-covered regions. This cooling is not fully captured by IFS during the evaluated period (Day et al., 2020).

Observation rank histograms are used to check if the forecasts cover the observed range at each station separately. This diagnostic includes all stations rather than focusing on the hottest or coldest temperatures in the verification domain as in a Q-Q plot. A flat histogram indicates that the distribution of forecasts and observed temperatures is similar. This is the case for the IFS forecasts over the summer in Fig. 6(b) while the tilted histogram for PGW reflects a systematic bias in the forecast. During winter, the IFS forecasts tend to be too cold at night time over mainland Europe (Sandu et al., 2020) while still not reaching extremely low temperatures in Northern Europe, as shown in Fig. 6. The overall negative
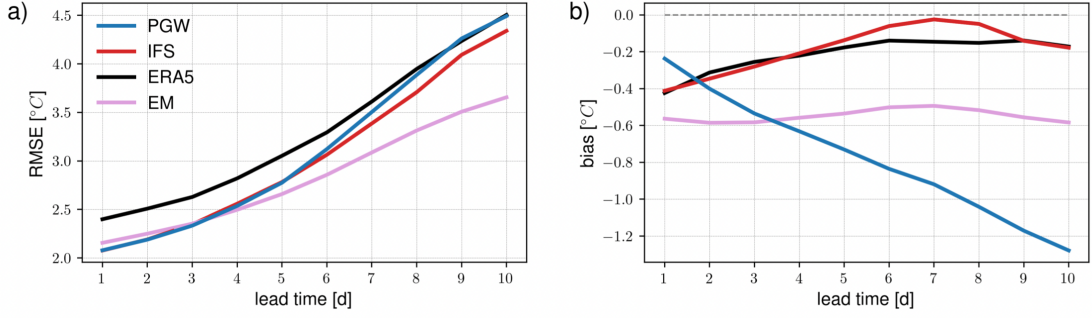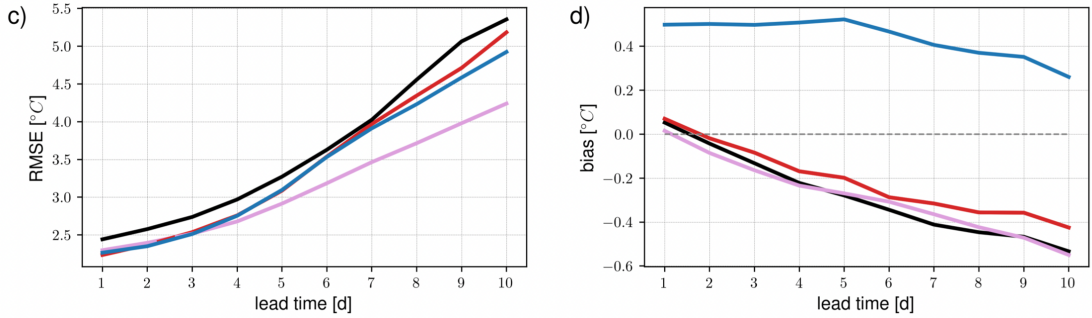
**Summer 2022**



**Winter 2022/2023**

Figure 5: Forecast performance for 2m temperature over Europe during Summer 2022 (top panels) and Winter 2022/2023 (bottom panels). (a,c) RMSE and (b,d) bias as a function of the forecast lead time. Summer 2022 forecasts are initialized at 12UTC (valid at midday) while Winter 2022/2023 forecasts are initialized at 00UTC (valid at midnight). The forecasts are verified against SYNOP observations.

bias leads to an over-populated first bin of the IFS histogram in Fig. 6(d) whilst PGW does not fully capture the lowest temperature at each station leading to underpopulated first bins of the histogram.

## 4.3 Forecasting weather events

The usefulness of a forecast is judged by its ability to predict weather events, often related to extremes. Here, the focus is on the forecast's ability to distinguish between an event and a non-event. Events are defined as 2m temperature exceeding a climate percentile. The climatology varies for each station and the climate percentiles are estimated based on the verification sample for the forecasts and the observations separately, to remove any bias in the forecast. We consider only low-temperature events for the winter period and high-temperature events for the summer period.

The relative operating characteristic (ROC) curve is a popular diagnostic tool in forecast verification. ROC curves plotting hit rate versus alarm rate of a high-temperature event in summer and a cold temperature event in winter are shown in Figs 7(a) and 7(c), respectively. Deterministic forecasts such as PGW and IFS forecasts have only one non-trivial point on the curve. This point is closer to the top-left corner of the plot for PGW than for IFS, indicating that PGW has better discrimination ability than IFS for the events under consideration. For the ensemble forecast, the ROC curve is built using one point for each probability issued by ENS using a standard 'trapezoidal' approach (Ben Bouallègue and Richardson, 2022). The ROC curve of a probabilistic forecast, represented here by the empty circles, covers a much wider area than the curve derived from a single forecast.

A standard measure for discrimination is the area under the ROC curve (AUC). In Figs 7(b,d), the AUC of 6-day ahead forecasts is plotted as a function of the percentile thresholds. The severity of the event increases as the climate percentiles get closer to 0% in winter and 100% in summer, indicating a rarer event under scrutiny. In general, AUC decreases when focusing on more intense/rare events as it becomes more difficult to predict such events with a deterministic
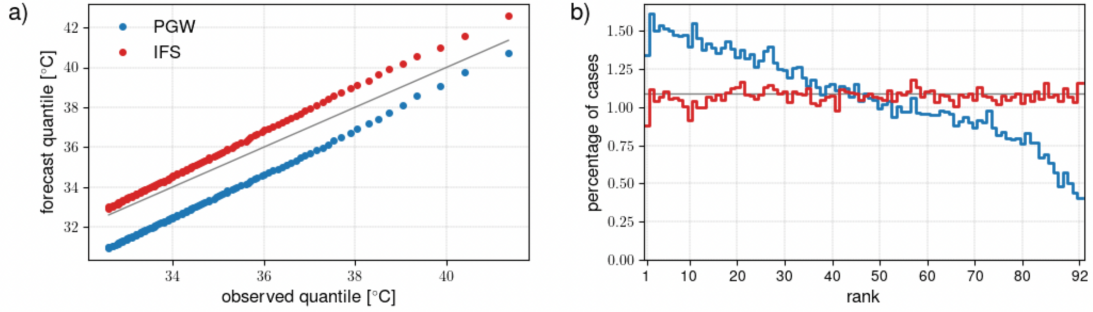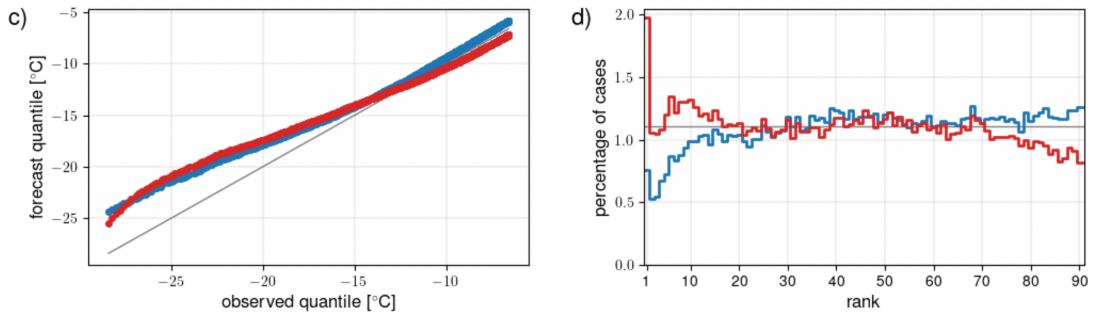
**Summer 2022**



**Winter 2022/2023**



Figure 6: Statistical consistency of 2m temperature over Europe during Summer 2022 (top panels) and Winter 2022/2023 (bottom panels) for IFS and PGW forecasts at day 6. (a,c) Q-Q plots showing a scatter plot of the empirical forecast quantiles versus the quantiles from the observation distribution at quantile levels 90%,90.1%, . . . ,99.9% for the summer period (a) and 0.1%,0.2%,..., 10.% for the winter period (c). (b,d) Observation rank histograms show the averaged number of forecasts in the bins defined by the sorted observations at each station. For all plots, perfect reliability is indicated by a grey line.

forecast. IFS and PGW have similar levels of performance in winter while PGW outperforms IFS in summer with differences statistically significant for percentiles between 75% and 95%, as estimated by block-bootstrapping.

## 4.4   Forecasting tropical cyclones

Tropical cyclones (TCs) are a prominent example of extreme weather that has a devastating impact and attracts considerable attention from the public and media. Moreover, TCs are characterized by large deviations from the mean state of the atmosphere and are thus generally challenging to forecast. Here, we focus on the year 2018 (as is done in Bi et al. (2023)), but note that IFS TC forecasts have substantially improved with more recent cycles (Forbes et al., 2021). We assess 2 key characteristics of TCs: their track position and their intensity (see the appendix for more details). In Fig. 8(a), the position error is measured as the distance between the TC position in the forecasts and the observations at a specific time. Larger errors are observed for PGW during the first day compared with IFS, but PGW has slightly lower errors for lead times greater than 2 days. This difference is partly explained by the fact that the propagation speed is generally too slow in the IFS (Chen et al., 2023a) but not in PGW (not shown). Overall, the differences in position error are small and not statistically significant between models.

Focusing now on TC intensity, Fig. 8(b) shows the mean absolute error for TC central pressure. Here we find that PGW clearly underestimates the intensity (*i.e.* the predicted pressure is too high). Both IFS and ERA5 perform better than PGW in terms of TC intensity error (except for at 0 day lead time). The large positive bias of PGW in the minimum core pressure results from too-weak gradients and too-weak maximum wind speed, while the IFS more closely resembles the analysis (not shown). The better performance of IFS compared with ERA5 is mainly explained by its higher resolution (9 km vs 28 km) but is also due to improvements in IFS through model development. From our investigations, PGW appears to have fewer TCs in its prediction than IFS. Further investigations are needed to explore the root causes for
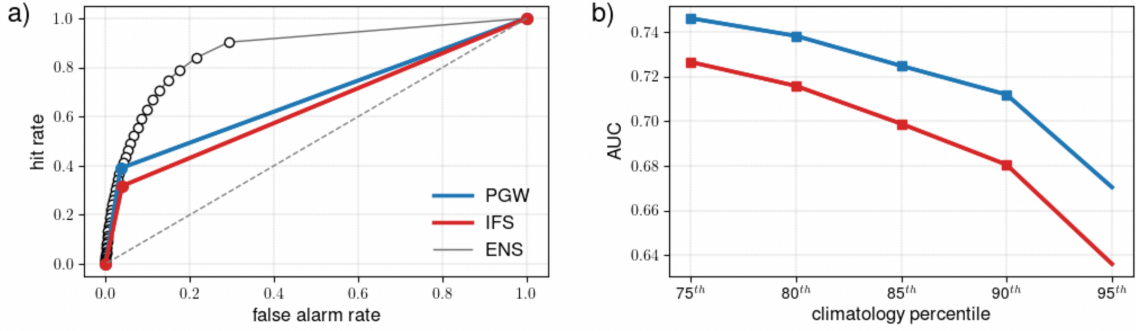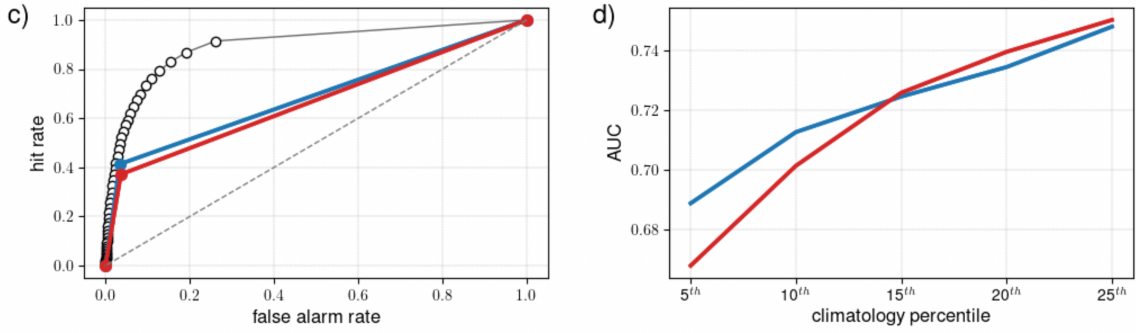
**Summer 2022**



**Winter 2022/2023**



Figure 7: Performance in forecasting 2m temperature events defined by climate thresholds for Summer 2022 (top panels) and Winter 2022/2023 (bottom panels) over Europe at day 6 lead time. (a,c) ROC curve (the closer to the top left corner the better) for an event defined as exceeding the 95% climate percentile in Summer (a) and below the 5% climate percentile in Winter (c). The diagonal dashed line is the zero-discrimination line. Results for ENS-derived probability forecast are also shown. (b,d) Discrimination ability as measured with the area under the ROC curve (AUC, the higher the better) and plotted as a function of the climate percentile used to define a weather event. A statistically significant difference between PGW and IFS results (as estimated by block-bootstrapping) is indicated by a square.

the reduced number of TCs in the ML-based forecast as well as to closely examine the TC structures and physical consistency between variables.

## 4.5 Predicting the forecast error

The day-to-day variability of the error is compared for PGW and IFS forecasts. We aim to identify common patterns in error growth and examine the sensitivity to predictability barriers. For this purpose, we analyze the so-called *predictability barrier plots* as described below. A more in-depth analysis would involve running different models from different initial conditions as in Magnusson et al. (2019) but this approach is out of scope for this paper.

Examples of predictability barrier plots for PWG and IFS are provided in Figs 9(a) and 9(b), respectively, focusing on daily scores of Z500 forecasts over Europe. In these plots, a transverse structure indicates rapid error growth leading to a poor forecast at all lead times: in that case, the forecast initialization might be the dominant predictability limiting factor. By contrast, a vertical structure indicates a weather situation difficult to predict for consecutive runs with different initialization, likely to be due to predictability barriers for that specific weather situation.

In general, we see a good agreement between PGW and IFS daily errors. This similarity is even more evident when plotting daily errors for a single lead time (here day 6) for the whole verification period. The correlation coefficient between the two time-series is 0.54. Strikingly, Fig. 9(c) shows the same 'bust' in forecasting the weather over Europe for 6[th] February. This flow-dependent nature of the error points towards the need for ensemble forecasting in a similar fashion for ML models like PWG as is common practice for NWP nowadays.
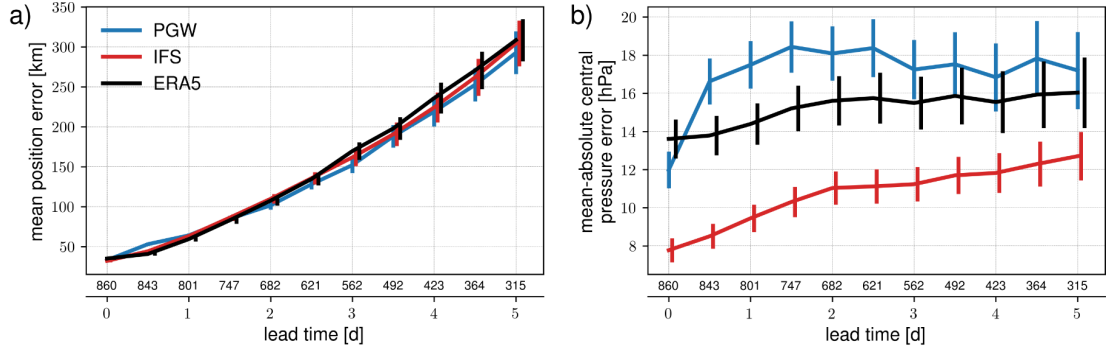
Figure 8: Tropical cyclone verification results: (a) mean position error and (b) mean absolute central pressure error as a function of the lead time for 2018. Forecasts are verified against the IBTrACS dataset and homogenized to have a consistent number of cases between models. For each lead time, the number of cases is displayed directly below the graphs. The vertical bars indicate the 2.5%-97.5% confidence intervals.

## 5    Summary and outlook

The results shown here highlight that machine learning (ML) models could have a promising future in numerical weather prediction. To explore the advantages and limitations of data-driven weather forecasts, we have run PanguWeather, an ML model trained on ERA5, initialized with the operational IFS analysis. The PWG forecasts are compared with the operational IFS forecasts to help shape our understanding of the characteristics of both the data-driven forecasts and their errors. Some of the most challenging weather phenomena are linked to rain, but our comparison does not include precipitation because the field is not present in PanguWeather.

Fundamentally, data-driven forecasts show good performance with the ML model being skillful for both upper-air variables (geo-potential height at 500hPa and temperature at 850hPa) when verified against operational analysis and for a surface variable (2m temperature) against observations. These conclusions are supported by further investigations including for other variables such as 10m wind speed. These results are not shown here because they are in line with the findings of this study.

We note, however, that there is room for further improvement in the implementation of data-driven weather prediction systems. Like all data-driven approaches presented so far, the ML model used in our study is not trained on the operational analysis and our experiments do not involve fine-tuning. Instead, ERA5 is used as the training dataset, which has so far been the cornerstone of any data-driven approach. For some aspects, ML models appear to directly inherit advantages and drawbacks from the numerical weather prediction system used to generate the training dataset. For example, the model resolution of the training data can in part explain the limitation in forecasting small-scale structures with PanguWeather. However, the forecast initial conditions also play a crucial role: starting a forecast from the operational IFS analysis rather than ERA5 analysis offers an advantage in skill also for the medium range. Moreover, the similarities in error growth of a data-driven forecast and a standard NWP forecast indicate similar sensitivities to chaos between ML-based and physically-based models.

The data-driven forecast appears smoother than the operational IFS forecast but the level of smoothness does not seem to increase with the forecast lead time, as we might expect when training toward RMSE. However, we observe a drift in bias almost linear with the forecast lead time. While this drift could be addressed when developing future ML models, statistical post-processing offers a means to correct systematic errors and achieve improved forecasts (Vannitsem et al., 2021). A deeper understanding of systematic errors could be achieved by performing conditional verification, for example, focusing on specific physical processes. Moreover, other diagnostic tools could be used to check for physical consistency based, for instance, on multivariate verification that accounts for the correlation between variables.

Good performance of data-driven forecasts is also observed in predicting some extreme events and confirmed by case studies. The results shown here focus first on events defined as climate threshold exceeded at a station location. The performance of ML-based models in forecasting TCs is under scrutiny too, as in (Bi et al., 2023). Preliminary investigations indicate that ML models could tend to predict fewer TCs compared to IFS, with tracks of similar quality, but IFS better captures their intensity and structure. Additional studies would help to demonstrate the value of data-driven forecasts as well as their strengths and weaknesses in supporting decision-making.
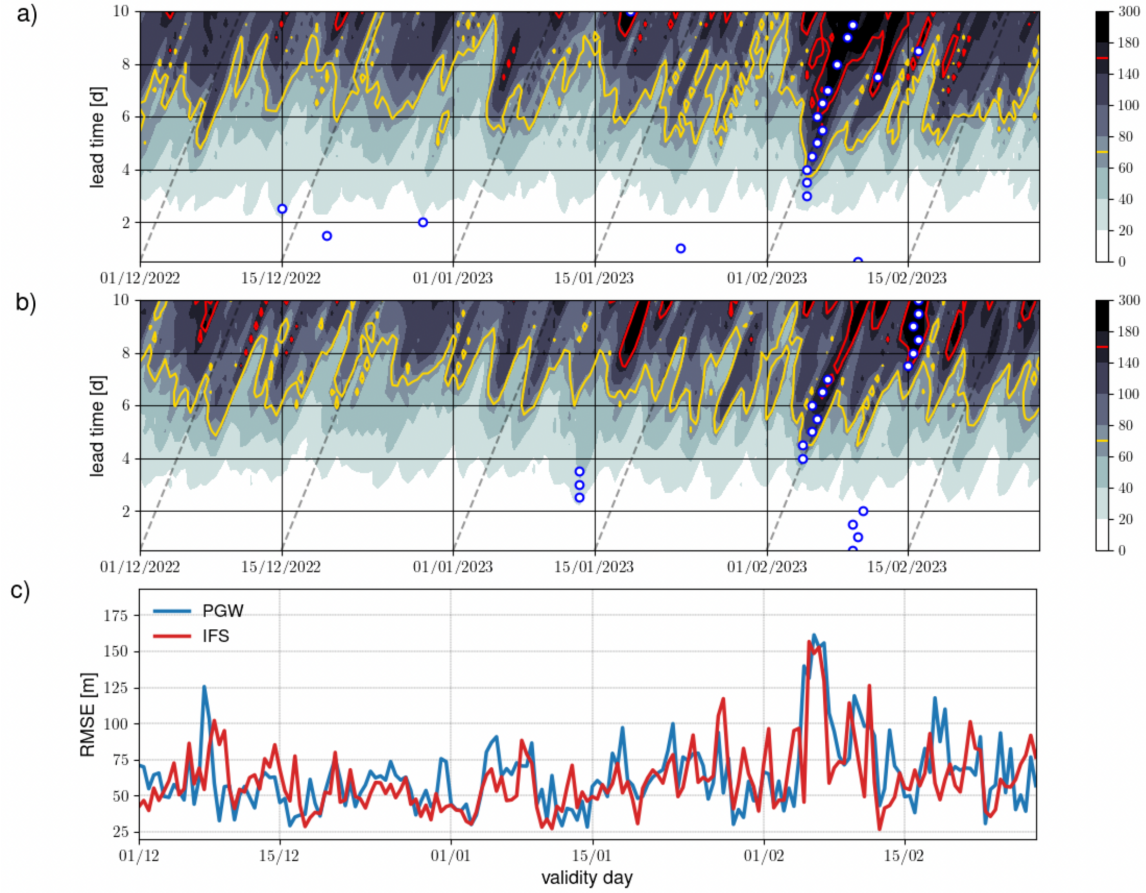
Figure 9: Predictability barrier plots showing daily RMSE for Z500 for lead times 1 to 10 days over Winter 2022/2023 for PWG (a) and IFS (b). A cross-section at day 6 of Figs (a) and (b) is provided in (c). In (a) and (b), the shade indicates the score value (in m), the vertical lines intercept scores for all forecasts valid on a given day, the transversal lines intercept scores for a given forecast run for all lead times, the yellow lines indicate the averaged score for a day 6 forecast, the red lines mark a large error, and the blue dots indicate the worst score over the period for each lead time.

Finally, in this work, we focused on deterministic forecasts but ensemble forecasts are key in providing uncertainty information for decision-making. A Monte-Carlo approach for uncertainty quantification has been tested starting ML-derived forecasts from perturbed initial conditions based on the ECMWF ensemble data assimilation and singular vector perturbations. The initial condition perturbation methodology is described in Lang et al. (2021b). The resulting ensemble forecast is showing promising results. As a future work, uncertainty in initial conditions will be complemented by mechanisms to account for model uncertainty (see e.g. Lang et al. (2021a)) in a data-driven weather prediction context.

This first assessment of a machine learning-based weather forecast in an operational-like context shows very promising results. The future role of ML models in the context of numerical weather prediction systems, and the ability of this approach to complement physical models remains to be explored. Operational centres should explore the strengths and weaknesses of these models as additional components of their forecasting systems: the ability to run forecasts at a much higher speed and much lower computational cost opens new horizons.

# References

Abbe, C., 1901: The physical basis of long-range weather forecasts. *Monthly Weather Review*, **29 (12)**, 551 – 561, doi:https://doi.org/10.1175/1520-0493(1901)29[551c:TPBOLW]2.0.CO;2.

Bauer, P., A. Thorpe, and G. Brunet, 2015: The quiet revolution of numerical weather prediction. *Nature*, **525**, 47–55.

Ben Bouallègue, Z., L. Magnusson, T. Haiden, and D. S. Richardson, 2019: Monitoring trends in ensemble forecast performance focusing on surface variables and high-impact events. *Quart. J. Roy. Meteor. Soc.*, **145 (721)**, 1741–1755, doi:10.1002/qj.3523.

Ben Bouallègue, Z., and D. S. Richardson, 2022: On the roc area of ensemble forecasts for rare events. *Weather and Forecasting*, **37 (5)**, 787 – 796, doi:10.1175/WAF-D-21-0195.1.

Bi, K., L. Xie, H. Zhang, X. Chen, X. Gu, and Q. Tian, 2023: Accurate medium-range global weather forecasting with 3D neural networks. *Nature*, doi:https://doi.org/10.1038/s41586-023-06185-3.

Bjerknes, V., 1904: Das Problem der Wettervorhersage, betrachtet vom Standpunkte der Mechanik und der Physik. *Meteorologische Zeitschrift*, **21**, 1–7.

Chen, J.-H., L. Zhou, L. Magnusson, R. McTaggart-Cowan, and M. Koehler, 2023a: Tropical cyclone forecasts in the dimosic project—medium-range forecast models with common initial conditions. *Earth and Space Science*, **10 (7)**, e2023EA002 821, doi:https://doi.org/10.1029/2023EA002821.

Chen, K., and Coauthors, 2023b: Fengwu: Pushing the skillful global medium-range weather forecast beyond 10 days lead. *arXiv preprint arXiv:2304.02948*.

Day, J. J., G. Arduini, I. Sandu, L. Magnusson, A. Beljaars, G. Balsamo, M. Rodwell, and D. Richardson, 2020: Measuring the impact of a new snow model using surface energy budget process relationships. *Journal of Advances in Modeling Earth Systems*, **12 (12)**, e2020MS002 144, doi:https://doi.org/10.1029/2020MS002144.

de Burgh-Day, C. O., and T. Leeuwenburg, 2023: Machine learning for numerical weather and climate modelling: a review. *EGUsphere*, **2023**, 1–48, doi:10.5194/egusphere-2023-350, URL https://egusphere.copernicus.org/preprints/2023/egusphere-2023-350/.

Dosovitskiy, A., and Coauthors, 2020: An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*.

Forbes, R., P. Laloyaux, and M. Rodwell, 2021: IFS upgrade improves moist physics and use of satellite observations. *ECMWF Newsletter*, **169**.

Geer, A. J., 2016: Significance of changes in medium-range forecast scores. *Tellus A: Dynamic Meteorology and Oceanography*, doi:10.3402/tellusa.v68.30229.

Hamill, T. M., and J. Juras, 2006: Measuring forecast skill: is it real or is it the varying climatology? *Quart. J. Roy. Meteor. Soc.*, **132**, 2905–2923, doi:doi.org/10.1256/qj.06.25.

Harvey, L. O., J. K. Hammond, C. Lusk, and E. Mross, 1992: The application of signal detection theory to weather forecasting behavior. *Mon. Wea. Rev.*, **120**, 863–883, doi:10.1175/1520-0493(1992)120<0863:TAOSDT>2.0.CO;2.

Hersbach, H., 2023: ERA5 reanalysis now available from 1940. *ECMWF Newsletter*, **175**.

Hersbach, H., and Coauthors, 2020: The ERA5 global reanalysis. *Quarterly Journal of the Royal Meteorological Society*, **146 (730)**, 1999–2049, doi:https://doi.org/10.1002/qj.3803.

Ingleby, B., 2015: Global assimilation of air temperature, humidity, wind and pressure from surface stations. *Quarterly Journal of the Royal Meteorological Society*, **141 (687)**, 504–517, doi:https://doi.org/10.1002/qj.2372.

Jolliffe, I. T., and D. B. Stephenson, 2011: *Forecast Verification: A Practitioner's Guide in Atmospheric Science, 2nd Edition*. 2nd Edn. Academic Press, New York, 627pp.

Keisler, R., 2022: Forecasting global weather with graph neural networks. *arXiv preprint arXiv:2202.07575*.

Knapp, K. R., H. J. Diamond, J. P. Kossin, M. C. Kruk, and C. J. I. Schreck, 2018: doi:https://doi.org/10.25921/82ty-9e16 [access date: July 2023].

Knapp, K. R., M. C. Kruk, D. H. Levinson, H. J. Diamond, and C. J. Neumann, 2010: The International Best Track Archive for Climate Stewardship (IBTrACS): Unifying Tropical Cyclone Data. *Bulletin of the American Meteorological Society*, **91 (3)**, 363 – 376, doi:https://doi.org/10.1175/2009BAMS2755.1.

Lam, R., and Coauthors, 2022: Graphcast: Learning skillful medium-range global weather forecasting. *arXiv preprint arXiv:2212.12794*.

Lang, S. T. K., S.-J. Lock, M. Leutbecher, P. Bechtold, and R. M. Forbes, 2021a: Revision of the stochastically perturbed parametrisations model uncertainty scheme in the integrated forecasting system. *Quarterly Journal of the Royal Meteorological Society*, **147 (735)**, 1364–1381, doi:https://doi.org/10.1002/qj.3978.

Lang, S. T. K., and Coauthors, 2021b: More accuracy with less precision. *Quarterly Journal of the Royal Meteorological Society*, **147 (741)**, 4358–4370, doi:https://doi.org/10.1002/qj.4181.

Leutbecher, M., and Z. Ben Bouallègue, 2020: On the probabilistic skill of dual-resolution ensemble forecasts. *Quarterly Journal of the Royal Meteorological Society*, **146 (727)**, 707–723, doi:https://doi.org/10.1002/qj.3704.

Leutbecher, M., and T. Palmer, 2008: Ensemble forecasting. *Journal of Computational Physics*, **227 (7)**, 3515–3539, doi:https://doi.org/10.1016/j.jcp.2007.02.014, URL https://www.sciencedirect.com/science/article/pii/S0021999107000812, predicting weather, climate and extreme events.

Magnusson, L., 2023: First exploration of forecasts for extreme weather cases with data-driven models at ECMWF. *ECMWF Newsl.*, **176**, doi:https://www.ecmwf.int/en/newsletter/176/news/.

Magnusson, L., J.-H. Chen, S.-J. Lin, L. Zhou, and X. Chen, 2019: Dependence on initial conditions versus model formulations for medium-range forecast error variations. *Quarterly Journal of the Royal Meteorological Society*, **145 (722)**, 2085–2100, doi:https://doi.org/10.1002/qj.3545.

Magnusson, L., and Coauthors, 2021: Tropical cyclone activities at ECMWF. *ECMWF Technical Memorandum*, **888**.

Mason, I., 1982: A model for assessment of weather forecasts. *Aust. Met. Mag.*, **30**, 291–303.

McGovern, A., R. Lagerquist, D. J. Gagne, G. E. Jergensen, K. L. Elmore, C. R. Homeyer, and T. Smith, 2019: Making the black box more transparent: Understanding the physical implications of machine learning. *Bulletin of the American Meteorological Society*, **100 (11)**, 2175 – 2199, doi:https://doi.org/10.1175/BAMS-D-18-0195.1.

Murphy, A. H., and R. L. Winkler, 1987: A general framework for forecast verification. *Mon. Wea. Rev.*, **115**, 1330–1338.

Pathak, J., and Coauthors, 2022: Fourcastnet: A global data-driven high-resolution weather model using adaptive fourier neural operators. *arXiv preprint arXiv:2202.11214*.

Sandu, I., and Coauthors, 2020: Addressing near-surface forecast biases: outcomes of the ECMWF project 'Understanding uncertainties in surface atmosphere exchange' (USURF). *ECMWF Technical Memorandum*, **875**.

Vannitsem, S., and Coauthors, 2021: Statistical postprocessing for weather forecasts: Review, challenges, and avenues in a big data world. *Bulletin of the American Meteorological Society*, **102 (3)**, E681–E699, doi:10.1175/BAMS-D-19-0308.1.

Wilks, D. S., 2006: *Statistical methods in the atmospheric sciences*. 2nd Edn. Academic Press, New York, 627pp.

## Appendix: verification process and scores definition

### Contextualising the forecast skill

We aim at a quantitative comparison of PGW and IFS forecasts. In general terms, forecast verification consists of measuring the relationship between a forecast and the corresponding observation[2] (Murphy and Winkler, 1987). For this purpose, one can carefully choose from a variety of metrics and diagnostics (see Wilks, 2006; Jolliffe and Stephenson, 2011). To go beyond the computation of generic scores, it is possible to investigate the properties of the joint distribution of forecasts and observations. The two main forecast attributes are forecast consistency (or calibration) and forecast discrimination ability. Note that these concepts hold also when dealing with probabilistic forecasts (which are not explored here).

Classical statistic tools involve the computation of summary statistics and scores. Summary statistics include the bias defined as the averaged difference between forecasts and observations, and forecast/observation activity defined as the standard deviation of the forecast/observation anomaly. Scores are metrics measuring the forecast skill such as the anomaly correlation shown in Fig. 1 or the forecast error such as the widely-used RMSE shown in Fig. 2. A formal definition of each of these metrics is provided in the Appendix. Here we formally define the following quantities:

- the forecast root mean squared error:

$$\sqrt{\overline{(f - o)^2}}, \tag{1}$$

- the forecast mean error (or bias):

$$\overline{f - o}, \tag{2}$$

- the forecast activity:

$$\sqrt{\overline{\left[(f - c) - \overline{f - c}\right]^2}}, \tag{3}$$

---

[2]the term *observation* is used in the broad sense of an assumed 'truth'. It can take the form of an analysis or an observation.

- the observation activity:

$$\sqrt{\overline{[(o-c)-\overline{o-c}]^2}},\qquad(4)$$

- the forecast anomaly correlation:

$$\frac{\overline{(f-c-\overline{f-c})(o-c-\overline{o-c})}}{\sqrt{\overline{(f-c-\overline{f-c})^2}}\sqrt{\overline{(o-c-\overline{o-c})^2}}}\qquad(5)$$

with $f$ the forecast, $o$ the observation, $c$ the climatology, and $\bar{\cdot}$ the averaging operator including a latitude weighting.

**Checking for statistical consistency**

Statistical consistency is tested regionally with quantile-quantile (Q-Q) plots and locally (at the station level) with observation rank histograms. For this exercise, we exclude stations situated at an altitude greater than 1000m to avoid focusing predominantly on representativeness issues rather than model characteristics. For Q-Q plots, quantiles are estimated from the whole verification sample (Europe, Summer 2022 or Winter 2022/2023) for both observations and forecasts, separately. We restrict our analysis to the warm tail of the distribution in the summer (quantile levels in the range $90\% - 99.9\%$) and to the cold tail of the distribution in the winter (quantile levels in the range $0.1\% - 10\%$).

As a complementary diagnostic tool, we suggest a new type of plot: the observation rank histogram (ORH). Inspired by the ensemble rank histogram used to assess the reliability of ensemble forecasts, ORH is built by ranking the observations from the smallest to the biggest for the whole verification period and individual forecasts, for each station separately. The rank of the forecast for each verification day is registered and populates the histogram. ORH assesses whether forecasts and observations are distributed similarly at a station level.

**Forecasting weather events**

Forecasting specific events is at the heart of many weather applications. The ability of a forecast to distinguish between the occurrence and non-occurrence of an event is called discrimination. Based on local climatology, event thresholds are defined as discussed above. Forecasts and observations are transformed into binary values with respect to a given threshold. A contingency table is populated for each of the dichotomous events. A contingency table is a $2 \times 2$ table where hits, misses, false alarms, and correct negatives are counted. From this table, it is possible to derive both the hit rate and the false alarm rate, the two components of the relative operating characteristic (ROC) curve. The area under the ROC curve (AUC) is a common measure of discrimination in weather forecast verification (Mason, 1982; Harvey et al., 1992).

Weather events are defined with the help of a local climatology that differs for each station. The same verification setting is used as in Ben Bouallègue et al. (2019) where an event is defined using a percentile of a climatology rather than a fixed absolute value. This approach tries to reflect that user-relevant thresholds are often associated with potential hazards and as such vary from place to place. For example, the 5% percentile of the local temperature climatology corresponds to very different absolute thresholds for say Helsinki and Madrid. Also using a climatology-based threshold allows us to avoid the pitfall of measuring varying climatology rather than actual skill (Hamill and Juras, 2006).

A different climatology is defined for (i) the observations and (ii) each forecast, with percentiles directly estimated from the verification sample. This so-called *eigen-climatology* approach corresponds to practically applying an in-sample local bias correction of the forecast as discussed in more detail in Ben Bouallègue et al. (2019). This step is important to disentangle discrimination from calibration attributes, because the latter can, in principle, be improved by post-processing. Only stations where measurements are available throughout the full verification period are considered for this exercise.

Finally, we recall that statistical significance is important when comparing competing forecasts (Geer, 2016). Here, we assess the chaotic variability of the scores with the use of (block)-bootstrapping. We randomly choose the verification days entering the verification dataset and compute scores for each forecast. Based on a 1000-member block bootstrap sample with blocks of 5 days, statistical significance to the 5% level is estimated.

**Forecasting tropical cyclones**

We also assess performance in forecasting tropical cyclones (TCs). TCs are tracked in forecasts from PGW, IFS, and ERA5 with the ECMWF operational TC tracker as described in Magnusson et al. (2021). Forecasts up to 5 days are verified here as results for longer lead times are unlikely to be statistically significant. As observations, we use

the International Best Track Archive for Climate Stewardship (IBTrACS) database (Knapp et al., 2010, 2018). The verification is based on TCs that are present in the observation database at the forecast initial time. The sample is homogenized to include the same cases for all 3 models. This homogenization results in a sample size of 860 cases at the analysis time, down to 315 cases at day 5[3]. An intensity threshold of 17m/s is applied to filter the observation dataset for TCs that reach tropical storm strength.

---

[3]If only the IFS was validated, the maximum number of cases would be 988 at the forecast initial time and 592 for 5-day forecasts.