

PeakWeather: MeteoSwiss Weather Station Measurements for Spatiotemporal Deep Learning

Daniele Zambon^{1*}

Michele Cattaneo^{2*}

Ivan Marisca¹

Jonas Bhend²

Daniele Nerini²

Cesare Alippi^{1,3}

¹ Università della Svizzera italiana, IDSIA, Lugano, Switzerland.
{daniele.zambon, ivan.marisca, cesare.alippi}@usi.ch

² Federal Office of Meteorology and Climatology MeteoSwiss, Zurich, Switzerland.
{michele.cattaneo, jonas.bhend, daniele.nerini}@meteoswiss.ch

³ Politecnico di Milano, Milan, Italy.

Abstract

Accurate weather forecasts are essential for supporting a wide range of activities and decision-making processes, as well as mitigating the impacts of adverse weather events. While traditional numerical weather prediction (NWP) remains the cornerstone of operational forecasting, machine learning is emerging as a powerful alternative for fast, flexible, and scalable predictions. We introduce PeakWeather, a high-quality dataset of surface weather observations collected every 10 minutes over more than 8 years from the ground stations of the Federal Office of Meteorology and Climatology MeteoSwiss’s measurement network. The dataset includes a diverse set of meteorological variables from 302 station locations distributed across Switzerland’s complex topography and is complemented with topographical indices derived from digital height models for context. Ensemble forecasts from the currently operational high-resolution NWP model are provided as a baseline forecast against which to evaluate new approaches. The dataset’s richness supports a broad spectrum of spatiotemporal tasks, including time series forecasting at various scales, graph structure learning, imputation, and virtual sensing. As such, PeakWeather serves as a real-world benchmark to advance both foundational machine learning research, meteorology, and sensor-based applications.

1 Introduction

Weather forecasts provide essential information for protecting lives and properties, as well as supporting decision-making in everyday activities. The quality of these forecasts has steadily improved over recent decades, driven by advances in numerical weather prediction (NWP) and high-performance computing (HPC) [7]. NWP involves using mathematical models that simulate the weather based on current conditions of atmosphere, land surfaces and oceanic conditions. Both deterministic forecasts and ensembles of simulations are produced, with ensemble forecasts offering a probabilistic view of future states and serving as a critical tool for risk-informed decision-making. However, the complexity of the models, combined with the demand for high-resolution predictions in space and time, makes state-of-the-art NWP computationally expensive—often requiring several hours of runtime on modern HPC systems.

*Equal contribution.

In parallel, the rise of deep learning (DL) has brought remarkable progress to a range of fields, including meteorology [21, 8, 25, 26]. Beyond breakthroughs in model architectures, this progress has been enabled by the availability of long-term, high-quality reanalysis datasets. Reanalyses combine modern NWP models with historical observational data—from ground stations, satellites, radar, ships, and aircraft—to produce a consistent and comprehensive estimate of past atmospheric states. DL models can learn to emulate the underlying physics captured in these massive datasets, enabling accurate and scalable data-driven weather forecasting. One of the key advantages of such models is their efficiency: once trained, they are significantly cheaper to run operationally than classical physics-based NWP. Nonetheless, these models still rely on NWP-generated analyses at inference time to initialize forecasts.

To overcome the need for NWP-based analysis and reanalysis data altogether, Direct Observation Prediction (DOP) methods have been proposed [23]. DOP approaches aim to learn directly from raw observational data, eliminating the need for NWP data during both training and inference. When successfully implemented, DOP models can offer faster, more cost-effective forecasting pipelines and avoid potential biases introduced by physics-based simulations. DOP has shown promising results in short-range weather forecasting, particularly for lead times ranging from minutes to a few hours [29, 33], a domain known as nowcasting. A typical nowcasting application is the extrapolation of radar image sequences to predict imminent precipitation. While DOP has proven effective in short time horizons, NWP systems generally maintain an edge in predictive accuracy beyond a few hours.

We introduce **PeakWeather**, a benchmarking dataset designed to advance deep learning research in spatiotemporal modeling. The dataset consists of high-resolution, ground-based meteorological measurements collected from SwissMetNet [6]—the network of automatic weather stations operated by the Swiss Federal Office of Meteorology and Climatology, MeteoSwiss. It includes eight key weather variables related to temperature, humidity, precipitation, pressure, sunshine, and wind, recorded every 10 minutes over more than eight years. These observations come from 302 meticulously maintained and quality-validated stations, distributed throughout Switzerland, a country characterized by complex topography, shaped by the Swiss Alps, the Central Plateau, and the Jura Mountains. In addition to meteorological variables, PeakWeather includes topographic descriptors derived from digital elevation models, capturing the unique characteristics of each station’s surrounding area.

To the best of our knowledge, this is the first openly available dataset that combines long-term, high-frequency ground station data with topographic context and state-of-the-art NWP baselines in complex terrain. The complexity of Switzerland’s topography, combined with seasonalities and time variance of the meteorological quantities, pose unique challenges for modeling spatial and temporal dependencies. These factors make PeakWeather an ideal testbed for advancing data-driven methods in a wide array of machine learning tasks, including—but not limited to—time series forecasting, virtual sensing, graph structure learning, and missing data imputation. To support benchmarking, the dataset also includes ensemble forecasts from a state-of-the-art, high-resolution NWP model that has been operational at MeteoSwiss since May 2024. These forecasts serve as a strong physics-based reference point, enabling direct comparisons with machine learning approaches and promoting research in hybrid and alternative forecasting paradigms.

Our contributions are summarized as follows:

- Release of PeakWeather dataset [4]. The dataset contains: (i) A collection of high-quality validated multivariate time series from 302 weather stations comprising eight meteorological variables with a temporal resolution of 10 minutes acquired from January 1st, 2017 to March 31st, 2025; (ii) Topographical descriptors derived from a digital height model describing the Swiss terrain over a 50-meter grid; (iii) Numerical weather predictions from a state-of-the-art high-resolution NWP model.
- A framework-agnostic, easy-to-use and well-documented Python library [5] for a seamless interaction with the dataset. The library enables users to easily load, align, and preprocess the data, providing them in ready-to-use tensor formats suitable for deep learning workflows.
- Evaluation of predictive models for wind forecasting—a task of high practical importance in domains such as renewable energy, aviation, and severe weather management. Indeed, practitioners and researchers can consider and elaborate on different meteorological variables in their research. A range of deep learning architectures, including recurrent neural networks and spatiotemporal graph neural networks, are contrasted on the prediction task against forecasts from a state-of-the-art NWP model. Reported results establish valuable baselines

for future research and demonstrate the promising potential of deep learning models as viable alternatives to NWP models in specific weather prediction tasks.

The PeakWeather dataset is released under the CC-BY-4.0 license and the associated library under the BSD-3-Clause. With PeakWeather, we aim at facilitating the use of MeteoSwiss’ data that, starting from Spring 2025, is gradually making available as Open Government Data (OGD) [3].

2 Related work

One major category of meteorological datasets consists of gridded reanalysis data produced by NWP models, which assimilate past observations to forecast the weather state. One of the main driving forces behind the recent success of data-driven models is ECMWF Reanalysis v5 (ERA5) [2]. This dataset contains a detailed estimate of the global state of the atmosphere, land surface and oceans since the 1950s, based on the reanalysis of ECMWF’s Integrated Forecasting System (IFS). Derived from ERA5, WeatherBench [27] and its recent successor WeatherBench2 [28] offer standardized benchmarking datasets and tools, to facilitate the evaluation and comparison of machine learning models for weather forecasting. For the task of statistical postprocessing [32], where the goal is to correct systematic errors in NWP forecasts, EUPPBench [14] provides a dataset with aligned IFS forecasts and direct observations for a variety of meteorological variables. So far, however, only coarse-grained NWP data is available in EUPPBench and temporal resolution is limited to 6 hours.

Unlike the datasets above, PeakWeather supports DOP, enabling the development of models that rely solely on past observations and static features, making the final model completely independent from NWP models. While datasets such as WeatherReal [20], Weather2k [34], Monash Weather [18], and StationBench [30] share this DOP focus, PeakWeather distinguishes itself through its unique combination of characteristics not available in other datasets: quality-controlled data with a high spatial density and 10-minute temporal resolution, spanning over eight years across Switzerland’s challenging topography, and physics-based forecasts from an operational NWP model for comparison.

3 PeakWeather dataset

We introduce PeakWeather, a publicly available dataset hosted on Hugging Face [4]. The dataset is organized into three core components: (i) high-resolution historical surface observations, (ii) static topographic features derived from elevation data, and (iii) ensemble forecasts from an operational NWP model. To facilitate easy access and streamlined preprocessing, we also release an open-source Python library [5] tailored to work seamlessly with the dataset.

Meteorological observations at stations Meteorological observations come from the SwissMetNet network, Switzerland’s reference ground-based measurement system for weather and climate monitoring operated by MeteoSwiss. The network in PeakWeather comprises 302 stations in total: 160 standard meteorological stations measuring various weather and climate parameters, complemented by 142 additional stations measuring precipitation. Stations are labeled as either meteorological stations or rain gauges, as shown in Figure 1a. The average distance to the 5 nearest neighbors is 18km across the 160 core meteorological stations, and 12km when all stations are included. The stations follow the World Meteorological Organization (WMO) standards, use consistent sensors, and are regularly maintained. The measurement data undergo automatic and manual quality controls. From the range of weather parameters available, we select air temperature, relative humidity, atmospheric pressure, sunshine duration, wind speed and direction, wind gusts, and the amount of precipitation, as summarized in Table 1; the table also reports the unique identifier (Short name) consistent with the SwissMetNet naming convention and OGD program [3]. Statistics about parameters measured by different stations are later reported in Table 2 and discussed at the end of the section, highlighting variability across stations.

Topographic descriptors Features of the terrain play a crucial role in weather and climate processes across multiple scales and by integrating detailed terrain information, prediction accuracy of data-driven models can be enhanced as demonstrated in Section 5. To this end, PeakWeather includes a high-resolution digital height model (DHM) with a spatial resolution of 50 meters, capturing the bare ground elevation without vegetation and buildings. The DHM is derived from the DHM25 [1]

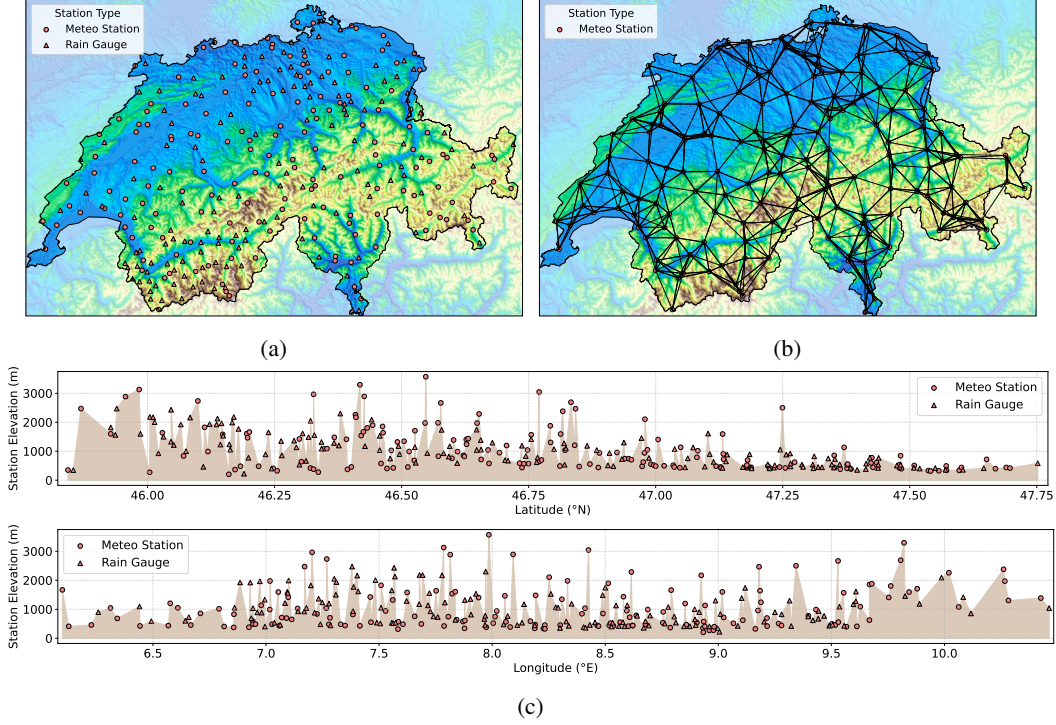


Figure 1: Visualizations of the placement of PeakWeather stations. Panel 1a) distribution of the stations across the Swiss territory; circles denote meteorological stations, while triangles denote rain gauges. Panel 1b) a graph obtained using the geographical distance to compute a similarity heuristic for the meteorological stations. Panel 1c) vertical coverage of the stations.

Table 1: Meteorological variables present in the dataset.

| Variable name | Short name | Description | Aggregation | Unit |
|----------------|------------|---|----------------|---------|
| temperature | tre200s0 | Air temperature 2m above the ground. | Current value. | Celsius |
| humidity | ure200s0 | Relative air humidity 2m above the ground. | Current value. | % |
| precipitation | rre150z0 | Precipitation. | 10min total. | mm |
| sunshine | sre000z0 | Sunshine duration. | 10min total. | min |
| pressure | presta0 | Atmospheric pressure at barometric altitude (QFE). | Current value. | hPa |
| wind_speed | fk1010z0 | Wind speed. | 10min mean | m/s |
| wind_gust | fk1010z1 | Wind gust peak of 1 second | 10min maximum. | m/s |
| wind_direction | dk1010z0 | Wind from direction, measured clockwise from north. | 10min mean. | degree |

product by the Swiss Federal Office of Topography swisstopo; a visualization is provided in Figure 1. In addition, we provide various derived topographical features that represent physical characteristics of the terrain and its surroundings. These features² are computed considering neighborhoods of 2 and 10 kilometers to cover phenomena at various scales and include the topographic position index (TPI) which describes whether a point is above or below the average elevation of its surroundings, the standard deviation (STD) of the elevation, the terrain aspect, the terrain slope, the west-east gradient of the terrain and the south-north gradient of the terrain.

Numerical weather predictions We include forecasts from ICON-CH1-EPS, one of the NWP models operational at MeteoSwiss since May 2024, as a strong physics-based baseline. ICON-CH1-EPS is a regional, high-resolution, and ensemble setup of the ICOSahedral Non-hydrostatic [ICON 35] model for Switzerland. It runs every 3 hours and produces an 11-member ensemble: one control run and ten perturbed members. Operating at $\sim 1\text{km}$ grid resolution, it generates forecasts up to 33 hours and the analysis step (lead time 0). PeakWeather includes hourly forecasts of the same variables and extracted at each station’s nearest model grid cell starting from May 14th, 2024.

²The implementation of these features has been open-sourced by MeteoSwiss <https://github.com/MeteoSwiss/topo-descriptors>.

Table 2: Number of stations per variable and station type, shown as “# in Jan 2017 / # in Mar 2025.” Missing data percentages are computed only for stations that recorded the variable.

| | Temp. | Humidity | Precip. | Sunshine | Pressure | Wind speed | Wind gusts | Wind dir. |
|----------------|---------|----------|---------|----------|----------|------------|------------|-----------|
| Rain gauges | 41/38 | 0/0 | 124/140 | 0/0 | 0/0 | 0/0 | 0/0 | 0/0 |
| Meteo stations | 148/149 | 148/149 | 138/139 | 103/131 | 137/139 | 149/151 | 150/151 | 150/151 |
| Missing values | 1.88% | 0.94% | 3.67% | 13.72% | 1.39% | 1.46% | 1.06% | 1.08% |

Practical considerations for data use One potential source of inconsistencies is sensor or station relocation, which—even in a carefully managed network like SwissMetNet—is sometimes necessary. To ensure transparency, PeakWeather also includes detailed metadata on all sensor relocations which, however, are generally minor: 95% of displacements are under 139 meters, with an average of 104 meters. Moreover, these relocations affect only a minority of stations, up to 17 for relative sunshine duration, and fewer than 10 for most other parameters. Although the stations and data are quality controlled, it should be noted that artifacts may still be present (e.g., due to frozen anemometers), and that the temperature sensors at rain gauges can be of inferior quality, as they are intended primarily to discriminate between solid and liquid precipitation.

Missing data Missing data in PeakWeather stems from both temporary measurement gaps and the fact that not all 302 SwissMetNet stations are equipped with the same set of sensors. Sensor configurations also evolve due to operational requirements. To give an accurate overview, Table 2 presents the number of stations recording each variable at the start and end of the dataset (January 2017 and March 2025), along with missing value statistics computed only for stations that have measured each variable at some point. Over the entire period, the overall proportion of missing data is approximately 3%. Missing values are represented as NaNs and can be handled using built-in utilities (e.g., last value carried forward or zero-fill). A binary availability mask is also provided by PeakWeather, allowing users to easily identify and handle invalid observations at each time step, station, and variable level.

Further details and visualizations are provided in Appendices A, B, and C.

4 Machine learning tasks on PeakWeather

The PeakWeather dataset enables a wide range of machine learning tasks. In this section, we outline several representative examples that can be addressed using the data. While not exhaustive, these tasks demonstrate the dataset’s versatility for both foundational research and practical applications in data-driven weather modeling.

We formalize the dataset as consisting of N multivariate time series, each collected from a different station. Let $\mathbf{x}_t^i \in \mathbb{R}^{D_x}$ denote the D_x -dimensional observation at time step t from station $i \in 1, \dots, N$; here, N denotes the total number of weather stations in the dataset and D_x the number of possible measured variables (channels) across all stations. Regular and synchronous sampling is assumed. However, as mentioned above, not all stations measure the same weather quantities, and not all sensors are available at all time steps (see Table 2). We account for such data (un)availability by appropriately padding the time series, thus maintaining a tabular representation. Accordingly, an auxiliary binary mask $\mathbf{m}_t^i \in \{0, 1\}^{D_x}$ is introduced at each time step t and station i to flag the data availability at the level of the time step, station, and channel. In particular, we set $\mathbf{m}_t^i[d] = 1$ if the d -th of the D_x channels of the i -th station at time step t is valid, and $\mathbf{m}_t^i[d] = 0$ otherwise. For conciseness, we denote by $\mathbf{X}_t \in \mathbb{R}^{N \times D_x}$ the stack of all observations $\{\mathbf{x}_t^i\}_{i=1}^N$ at time step t , by $\mathbf{X}_{t:t+T} \in \mathbb{R}^{T \times N \times D_x}$ the observations from time step in $\{t, t+1, \dots, t+T-1\}$ and by $\mathbf{X}_{<t}$ all observations up to time step t (excluded); similarly, $\mathbf{M}_t \in \mathbb{R}^{N \times D_x}$ denotes the mask across all stations at time step t . The PeakWeather library provides functionalities to generate such a mask.

As the time series correlates among themselves, we model the data-generating process as a time-invariant spatiotemporal stochastic process such that

$$\mathbf{x}_t^i = p^i(\mathbf{x}_t^i | \mathbf{X}_{<t}, \mathbf{U}_{\leq t}, \mathbf{V}), \quad \forall i = 1, \dots, N, \quad (1)$$

where $\mathbf{V} \in \mathbb{R}^{N \times d_v}$ and $\mathbf{U}_t \in \mathbb{R}^{N \times d_u}$ represent, respectively, static and time-dependent exogenous variables describing additional information related to the measured variables \mathbf{x}_t^i , such as station-specific characteristics, seasonalities, as well as the topographic descriptors from the DHM (see

Section 3). Note that the process of Equation 1 generating all time series is not necessarily the same, i.e., $p^i \neq p^j$ if $i \neq j$, while the assumption of time invariance remains valid.

Relational inductive biases for time series processing Time series processing can be significantly improved by exploiting the underlying (functional) dependencies among series acquired from spatially distributed stations. Rather than modeling each sequence in isolation, one can incorporate relational information as an inductive bias, allowing models to condition predictions on related time series and reduce the risk of overfitting to local, spurious patterns. These dependencies can be encoded via a graph structure represented by an adjacency matrix $\mathbf{A} \in \{0, 1\}^{N \times N}$, where each node corresponds to a station and non-zero entries denote pairwise relationships. When needed, \mathbf{A} may be enriched with edge weights or attributes reflecting the strength or type of interaction (e.g., distance, elevation difference, empirical correlation). This setup is at the core of graph deep learning approaches for time series processing [13], where message-passing operators like graph neural networks (GNN) are exploited to exchange information across the graph topology. In the context of meteorological data, such relations are often assumed to be time-invariant. However, changes to the sensor network—such as the addition, removal, or relocation of stations—may require considering a dynamic or evolving graph. While spatial proximity naturally induces correlation (i.e., signal smoothness), this assumption breaks down at larger scales or in complex terrains. In alpine regions, for instance, valleys, ridges, and altitude gradients introduce significant variability, making long-range dependencies non-trivial to infer. Particularly in datasets like PeakWeather, this motivates considering relations learned data-driven and tailored to the given task, as discussed in the next experimental section.

Time series forecasting This task consists of predicting future values of meteorological variables of interest (e.g., temperature, wind speed) at individual stations, possibly all, by using historical observations and, whenever available, exogenous variables affecting the system. Denote by $\mathbf{y}_t^i \in \mathbb{R}^{D_y}$ the vector of target variables and formulate the task as that of learning a model p_θ to predict \mathbf{y}_{t+h}^i as

$$\hat{\mathbf{y}}_{t+h}^i \sim p_\theta(\mathbf{y}_{t+h}^i | \mathbf{X}_{t-W:t}, \mathbf{M}_{t-W:t}, \mathbf{U}_{t-W:t+h+1}, \mathbf{V}) \approx p^i(\mathbf{y}_{t+h}^i | \mathbf{X}_{<t}, \mathbf{U}_{\leq t+h}, \mathbf{V}), \quad (2)$$

for a set of lead times $h = 0, \dots, H - 1$ and stations i of interest; $W \geq 1$ denotes the time window length used to make the predictions. This can be performed at multiple temporal resolutions and forecast horizons, and can leverage both spatial and temporal dependencies to predict weather conditions across multiple stations simultaneously. This includes using graph-based models to capture correlations between locations that condition predictions on a graph \mathbf{A} .

Wind forecasting is the application considered in the experiments below (Section 5). Accurate wind forecasts provide useful information for diverse applications such as anticipating extreme weather and optimizing renewable energy production. Forecasting wind in complex terrain, however, is challenging due to the dominant influence of local conditions. Exposure to prevailing large-scale winds and related effects, such as channeling and blocking, affects wind speeds at individual locations. In addition, distinct phenomena such as mountain valley wind systems dominate the temporal evolution.

Missing data imputation and virtual sensing Missing data imputation addresses temporal gaps caused by missing or erroneous sensor readings, aiming to reconstruct the time series at observed locations. In contrast, virtual sensing (or spatial interpolation) focuses on estimating weather variables at locations without prior measurements using spatial patterns from nearby stations. These tasks can involve conditioning on spatially and temporally related observations (possibly from future time steps too), topographical information, and other related variables at the stations to improve accuracy. Despite the high reliability of SwissMetNet, missing observations remain inevitable (Table 2). Moreover, the continuous updates to the SwissMetNet network, where new stations are added, some are reconfigured, and others are removed, motivate considering inductive learning methods, where model predictions can extend to new, unseen stations based on their spatial and contextual characteristics.

Graph structure learning As demonstrated in Section 5, better predictions can be issued by leveraging relational information among the time series. However, it is not always evident which station/variable would benefit from exchanging information. While it is common to extract a graph of binary relations from, e.g., physical proximity or correlations among time series, graph structure learning emerged as a task where relations can be treated as latent variables that are learned data-driven, optimizing a prediction task of interest. Given the complex Swiss terrain and the specificity

of each station’s location, performing graph structure learning in PeakWeather appears sound and insightful.

5 Experiments: a use case in wind forecasting

In this section, we address a wind forecasting task, selected here as a use-case spatiotemporal modeling task to showcase the potential of PeakWeather. In particular, we show how the historical measurements from the stations and the topographical features can be integrated to yield accurate predictions, as compared to the reference NWP model predictions.

Wind forecasting setup We consider the task of predicting wind direction and speed for the next 24 hours. Specifically, we consider hourly data and set the wind horizontal velocity vector as target $\mathbf{y}_t^i \in \mathbb{R}^2$. The components of \mathbf{y}_t^i , also called u and v components, are constructed from variables `wind_speed` and `wind_direction` and indicate the eastward and northward wind components. For this task, we consider data from the 160 meteorological stations and with static attributes to encode the latitude, longitude, station height, and other topographic features. Notably, resampling from 10-minute resolution to 1-hour (while accounting for variable-specific aggregation strategies), u and v component extraction, and weather station filtering can all be carried out via the utilities in the PeakWeather dataset library. The considered models described in the following sections are trained on the historical data until December 31st, 2023 (i.e., seven years), tested on the last year of available data, and validated on the remaining months. A graph similar to that in Figure 1b is constructed by connecting stations based on their haversine distance. Further details on the training procedure, model selection, and graph construction are provided in Appendix D.

Evaluation methodology We train models minimizing the Energy Score (ES) [17]—an extension of the Continuous Ranked Probability Score (CRPS) for multivariate predictions—which evaluates the quality of probabilistic forecasts accounting for both the calibration and sharpness of predicted distributions; the ES is estimated and optimized via Monte Carlo sampling. Models are then evaluated on the test set using both probabilistic and point prediction metrics. Specifically, we report the ES for the bivariate wind velocity, and the mean absolute error on the same vector (MAE-vel) to assess point-forecast accuracy. In addition, we include the MAE for the wind speed (MAE-speed) and for the wind direction (MAE-dir), the latter computed as the mean angular distance between the target and predicted wind velocity. All metrics are reported across several lead times up to 24 hours. MAE, MAE-speed, and MAE-dir are evaluated considering the median of the wind velocity, speed and direction, respectively, over 100 Monte Carlo samples (11 samples for ICON). Missing observations are appropriately masked during evaluation.

5.1 Forecasting models and baselines

We consider two main deep learning forecasting architectures, recurrent neural networks (RNNs) and spatiotemporal graph neural networks (STGNNs) [19], which we contrast with two persistence baseline models (PM) and the NWP model provided within PeakWeather.

The **RNN** models process each station independently from the others. It is a multilayer RNN equipped with gated recurrent units [10]. The architectures feature an encoder block to aggregate the input \mathbf{x}_{t-w}^i and exogenous variables \mathbf{u}_{t-w}^i and \mathbf{v}^i of station i for each input time step $t - w$, before feeding them to the RNN, and a decoder responsible for making forecasts at the desired lead times; both encoder and decoder are implemented as multilayer perceptrons. The **STGNN** models implement a time-then-space architecture [13], following the same RNN architecture above, but composing GNN layers on top of the RNN, before issuing predictions via the decoder. For both RNN and STGNN, the model outputs bidimensional vectors $\hat{\mathbf{y}}_{t+h}^i$ for each lead time h as point predictions. Instead, for probabilistic forecasts, for each lead time h independently, a bivariate Gaussian distribution is learned via the reparametrization trick and conditioned on the model inputs (i.e., $\mathbf{x}_{t-W:t}^i$, $\mathbf{u}_{t-W:t}^i$, and \mathbf{v}^i). Optionally, node-level learnable embeddings $\mathbf{e}^i \in \mathbb{R}^{D_e}$ are learned end-to-end alongside the rest of the model parameters and are used to condition the model encoder and decoder outputs to account for station-specific dynamics [12]; we discuss them in Section 5.2. Further details on the model architectures follow in the supplementary material.

Two versions of the PM are implemented: one particularly relevant for short-term, **PM-st**, and the other, **PM-day**, intended to replicate the daily patterns of the target variables. PM-st is a baseline model that predicts $\hat{\mathbf{y}}_{t+h}^i = \mathbf{y}_{t-1}^i$ as the last available observation of the target for each lead time $h = 0, \dots, H-1$, a valid assumption for very short-term forecasts. Conversely, the output of PM-day is $\hat{\mathbf{y}}_{t+h}^i = \mathbf{y}_{t+h-24}^i$, recalling that in this specific problem instance we are considering hourly data and forecasting horizon of $H = 24$. While less evident than with other weather quantities, the wind can have a diurnal seasonality, especially in mountainous valleys. For probabilistic forecasts, independent Gaussian distributions $\mathcal{N}(\hat{\mathbf{y}}_{t+h}^i[d], \sigma_{t-W:t}^i)$ are produced for every d, h, t and i , with $\sigma_{t-W:t}^i$ being the sampling standard deviation estimated on the time window $\mathbf{y}_{t-W:t}^i[d]$.

Finally, the **ICON** numerical predictions (see Section 3) provided within PeakWeather are compared with the above models; as ICON has been operational since May 2024 issuing predictions every 3 hours, an additional specific subset of the test set (NWP test set) will be considered for comparing with this model.

Sample 24-hour forecasts of trained models and ICON are shown in Appendix E.

5.2 Results

We present three sets of experiments to showcase how PeakWeather can be used and why it is a valuable resource for the deep learning and spatiotemporal modeling research communities. **(i)** We analyze the impact of incorporating station-specific information—a key feature of PeakWeather—highlighting the presence of local effects. **(ii)** We assess the wind prediction accuracy of different models, showcasing the competitiveness of graph deep learning models, like STGNNs, for spatiotemporal modeling, and providing baseline prediction accuracies for future research. Finally, **(iii)** we compare the best-performing models with forecasts from the ICON numerical weather prediction (NWP) model, demonstrating the challenging nature of the task.

Relevance of station-specific features Table 3 presents the test performance for RNN and STGNN models trained with and without station-specific attributes: static vectors \mathbf{v}^i of station descriptors and learnable node embeddings \mathbf{e}^i . For RNNs, incorporating either \mathbf{v}^i , \mathbf{e}^i , or both leads to substantial improvement in the predictive accuracy, highlighting the importance of station-level information in the absence of built-in spatial modeling. For STGNNs, however, the inclusion of station-specific features results in only marginal performance differences. This suggests that STGNNs already capture key spatial dependencies through message passing, making explicit station embeddings less critical. Nevertheless, STGNNs consistently outperform RNNs overall, confirming their advantage in spatiotemporal modeling.

Table 3: Assessment of the impact of station-specific features on the prediction performance. Results are aggregated over five runs with different random seeds and are reported as mean \pm standard deviation.

| Model | MAE-vel | ES |
|-----------------------------------|-------------------|-------------------|
| RNN | 1.173 \pm 0.005 | 1.334 \pm 0.008 |
| RNN(\mathbf{v}) | 1.134 \pm 0.006 | 1.281 \pm 0.007 |
| RNN(\mathbf{e}) | 1.143 \pm 0.006 | 1.292 \pm 0.009 |
| RNN(\mathbf{e}, \mathbf{v}) | 1.133 \pm 0.002 | 1.279 \pm 0.002 |
| STGNN | 1.096 \pm 0.001 | 1.241 \pm 0.004 |
| STGNN(\mathbf{v}) | 1.088 \pm 0.004 | 1.227 \pm 0.004 |
| STGNN(\mathbf{e}) | 1.094 \pm 0.003 | 1.233 \pm 0.001 |
| STGNN(\mathbf{e}, \mathbf{v}) | 1.094 \pm 0.006 | 1.236 \pm 0.011 |

Forecasting Accuracy The top part of Figure 2 shows forecasting performance across lead times for all models. All trainable models include static station attributes; among them, RNN and STGNN also utilize learnable node embeddings, while RNN-noemb is identical to RNN except for the omitted node embedding \mathbf{e}^i and discarded static exogenous \mathbf{v}^i . PM-st baseline provides a strong reference, particularly in the short term. As the horizon increases, PM-day marginally outperforms PM-st, although the difference is minor within the considered 24 hours; as expected, PM-st and PM-day coincide at 24h predictions. All deep learning models substantially outperform the PM baselines across all horizons, both in terms of MAE and probabilistic accuracy, measured by the ES metric. Among them, the STGNN model achieves the best overall performance, underscoring the value of incorporating spatial structure. As expected, performance degrades with increasing forecast horizon. This trend holds for all models, with the exception of the persistence baselines: PM-day effectively remains fixed at 24h forecasting horizon, while PM-st benefits from daily wind patterns. Notably, wind direction is consistently harder to predict than wind speed as evidenced by smaller improvements over the PM baselines.

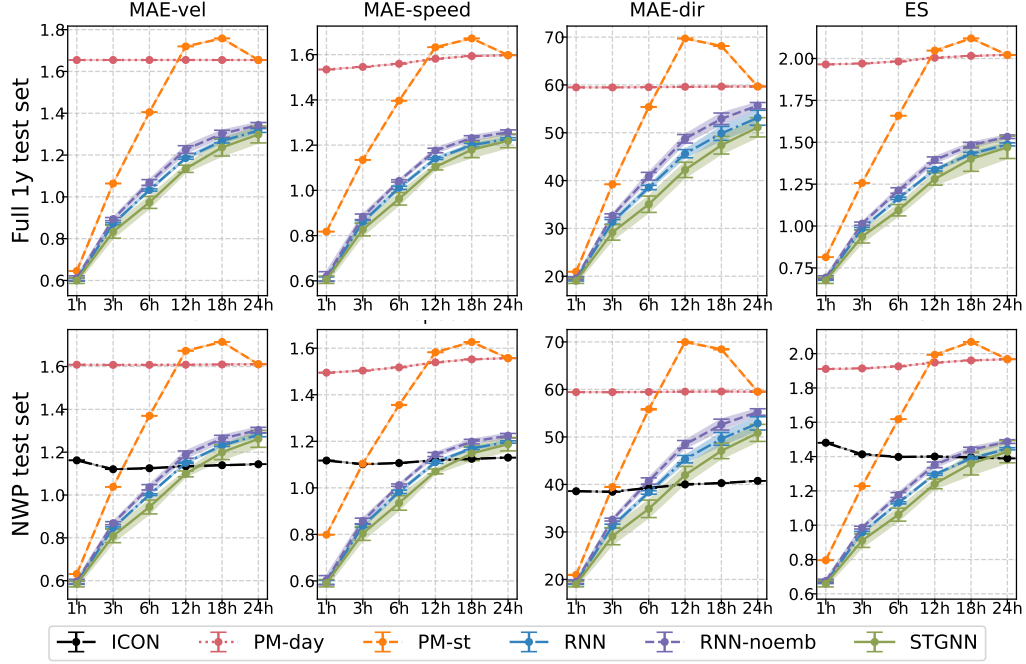


Figure 2: Performance comparison of deep learning models with ICON and persistence model baselines for wind forecasting. Results are averaged over five runs with different random seeds, except for ICON. As the persistence models require no training, their variability arises solely from Monte Carlo sampling during metric evaluation. Shaded areas indicate ± 3 standard deviations.

Comparison with ICON NWP forecasts Figure 2 also includes model performance on the NWP test set, the subset of the full 1-year test set aligned with ICON forecasts. First, we note that the forecasting metrics on this subset are consistent with those observed on the full test set. Then, we observe that ICON demonstrates relatively stable performance across all lead times, including short-term forecasts. However, its accuracy does not improve significantly at near-term horizons; this is likely due to spatial mismatches between the gridded NWP output and the fine-scale conditions at station locations and consequent challenges in assimilating such local information in the model. Comparing ICON to the deep learning models reveals that STGNN delivers more accurate predictions in the nowcasting regime, and up to around 12 hours ahead. For longer horizons, and particularly for wind direction, ICON outperforms the data-driven approaches, suggesting that physics-based models still hold an advantage in capturing larger-scale dynamics over extended periods.

These findings align with our earlier discussion: deep learning models excel at short-term forecasts, offering efficient and accurate predictions, while NWP models currently maintain superior performance for longer horizons. This highlights an important frontier for machine learning research in spatiotemporal modeling.

6 Conclusion

The paper introduces PeakWeather, a high-resolution dataset of validated ground-based weather measurements collected from the SwissMetNet network operated by MeteoSwiss. Designed to support and accelerate research in spatiotemporal deep learning, the dataset provides dense and diverse weather observations from 302 stations distributed across Switzerland. The dataset includes eight meteorological variables sampled every 10 minutes for more than eight years. Notably, the data is enriched with topographical descriptors to account for the complex Swiss terrain and with ensemble forecasts from a state-of-the-art NWP model currently operational at MeteoSwiss for benchmarking.

Beyond its scale and quality, PeakWeather distinguishes itself by supporting a wide range of machine learning tasks. These include time series forecasting, virtual sensing (spatial interpolation), graph structure learning, and missing data imputation. PeakWeather data is provided alongside a framework-

agnostic library that offers data in ready-to-use tensor formats for seamless integration into modern deep learning workflows. Experiments on wind forecasting conclude the paper, showcasing the timeliness of PeakWeather, its relevant features for spatiotemporal modeling and reinforcing the emerging role of deep learning as a complementary tool to traditional NWP systems. We are confident that PeakWeather will foster both fundamental machine learning research as well as advancing data-driven approaches in meteorological modeling.

Acknowledgments

This work was supported by the Swiss National Science Foundation project FNS 204061: *HORD GNN: Higher-Order Relations and Dynamics in Graph Neural Networks*.

A Details about the NWP forecasting model ICON-CH1-EPS

The numerical weather prediction (NWP) baseline from the MeteoSwiss’ ICON-CH1-EPS consists of hourly forecasts at stations. For each forecast initialization every 3 hours, forecasts for 11 ensemble members are available. The first ensemble member is the control run, which has been initialized and run with unperturbed boundary conditions (i.e., the best estimate of the state of the atmosphere). The remaining 10 ensemble members use randomly perturbed boundary conditions and stochastically perturbed physics. This also implies that there is no connection between realizations of successive initializations other than for the control member and perturbed ensemble members should be treated as statistically exchangeable.

The forecasts cover the time span from forecast initialization to 33 hours in the future. For the aggregated quantities precipitation, sunshine, and wind_gust, the analysis timestep (i.e., forecast lead time 0) is missing.

The ICON-CH1-EPS forecasts at stations have been extracted from the values at the nearest grid point of the regular rotated latitude-longitude grid used at MeteoSwiss for processing of high-resolution NWP data; the distance between each station and its corresponding grid point is less than 1 km. To account for the difference in model topography and altitude, temperature forecasts from ICON-CH1-EPS are altitude-corrected. As is practice at MeteoSwiss, we use a constant lapse-rate of $0.06K/m$ for altitude correction. Relative humidity is computed from temperature and dew-point temperature via the ratio of saturated vapour pressure of water derived using Eq. 10 from [9]. For the computation of relative humidity, neither temperature nor dew-point temperature is altitude corrected.

B Visualization of topographic descriptors

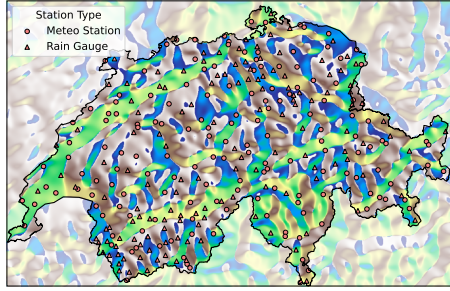
Figure 3 shows the topographic features included in PeakWeather, which are computed at both 2 km and 10 km spatial scales. The Aspect (Figures 3a and 3b) is defined as the azimuth (in degrees from North) of the steepest downslope direction. The Standard Deviation (STD) (Figures 3c and 3d) is the standard deviation of the surrounding elevation. The Topographic Position Index (TPI) (Figures 3e and 3f) describes whether a point is above or below the average elevation of its surroundings. The Slope (Figures 3g and 3h) is defined as the magnitude of the gradient vector of elevation. The West-East and South-North Derivatives (Figures 3i, 3j, 3k, and 3l) are the partial derivatives of elevation in the west-east and north-south directions, respectively. Further information and implementation details can be found at the following link: <https://github.com/MeteoSwiss/topo-descriptors>.

C Variables availability over time

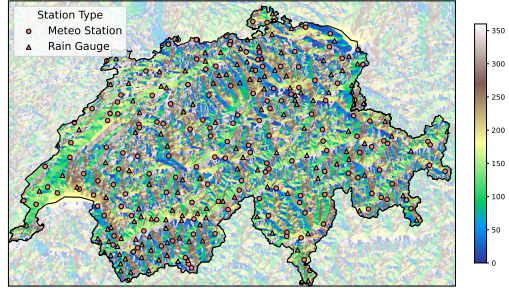
Figure 4 displays the number of stations per month that provide valid measurements for each of the 8 meteorological variables. The plot spans the entire dataset, from January 2017 to March 2025.

D Additional experimental details

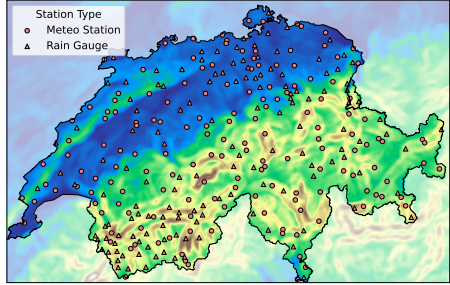
Model architecture The RNN and STGNN models have a 2-layer MLP as input encoder and 2-layer temporal processing via gated recurrent units with a hidden dimension of 128. The model



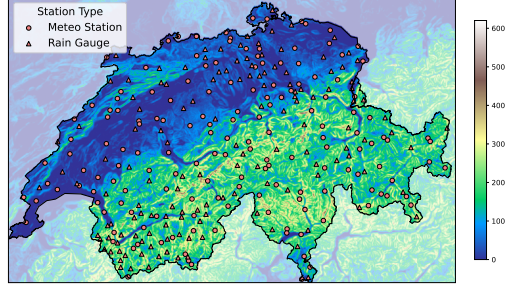
(a) Aspect 10km



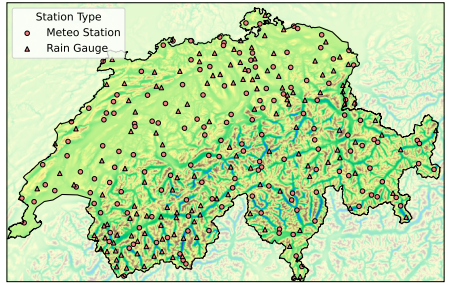
(b) Aspect 2km



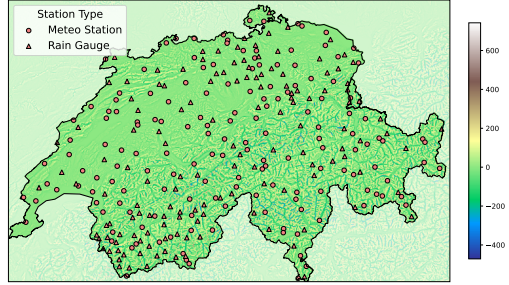
(c) STD 10km



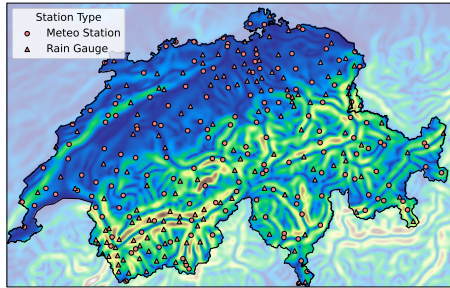
(d) STD 2km



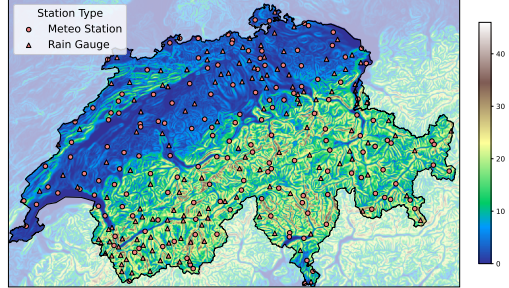
(e) TPI 10km



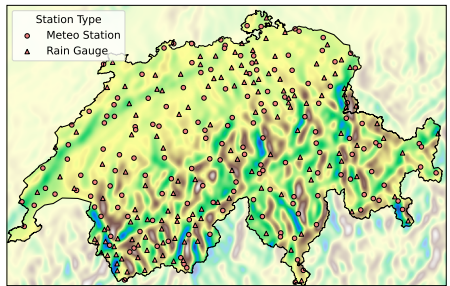
(f) TPI 2km



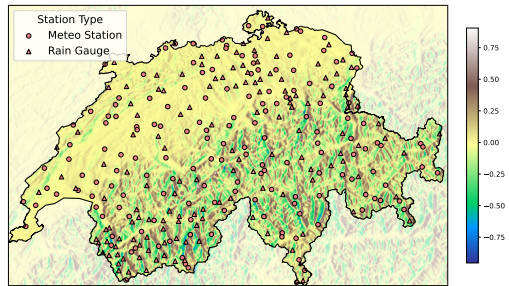
(g) Slope 10km



(h) Slope 2km



(i) WE-derivative 10km



(j) WE-derivative 2km

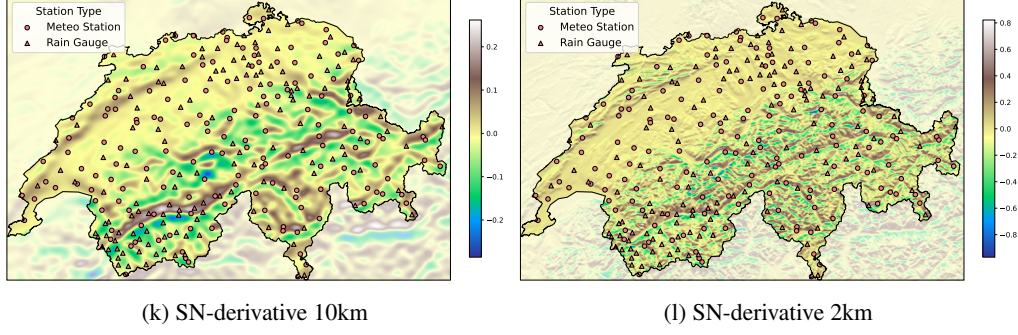


Figure 3: Visualization of the topographic descriptors.

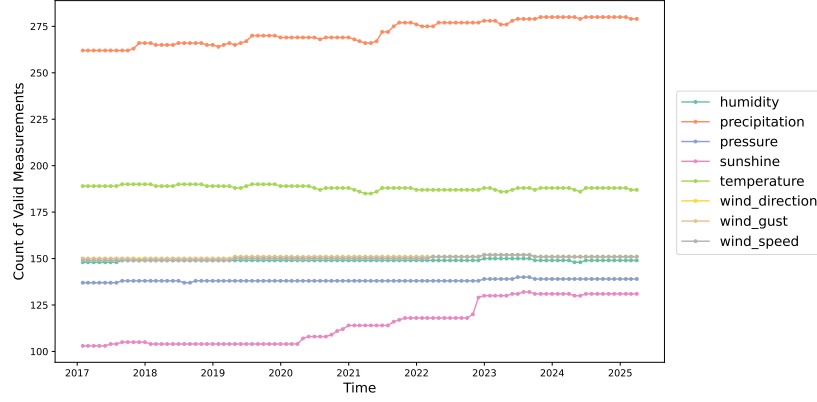


Figure 4: Stations with valid measurements of the meteorological variables included in PeakWeather.

predictions at each lead time and station are Gaussian distributions learned via the reparametrization trick, with mean and covariance conditioned on the inputs, lead time and station. The exponential linear unit is used as activation function, and 32-dimensional node embeddings are added as inputs to the encoder and readout layer, unless specified otherwise. STGNNs implement 6 message passing layers based on diffusion convolution from [22]. Model hyperparameters have been selected as those that yielded minimum MAE on the validation set.

Graph construction The graph used in the experiments is constructed based on pairwise geographical distances between stations. The Gaussian kernel with width equal to the standard deviation of the station distances is applied to compute station similarity scores. To enforce sparsity in the graph, for each node, edges are included only for stations with a similarity above a predefined threshold of 0.7, and if the stations are among the $k = 8$ nearest neighbors.

Model training and evaluation Model predictions are made for a 24-hour horizon from an input window of 24 time steps (1 day). All weather variables, encodings of the hour of the day and the day of the year, the binary mask of available observations, and all provided topographic descriptors are considered as input features to the model. Models are trained on the first 7 years of data by optimizing the Energy Score (ES) [17] at every station and every lead time (up to 24 hours ahead), validated on three months and tested on the last year, appropriately masking missing observations. The Energy Score is here estimated via Monte Carlo

$$\text{ES} \left(\left\{ \hat{\mathbf{y}}_{t+h}^{i,(m)} \right\}, \mathbf{y}_{t+h}^i \right) = \frac{\sum_{m=1}^M \left\| \hat{\mathbf{y}}_{t+h}^{i,(m)} - \mathbf{y}_{t+h}^i \right\|_2}{M} - \frac{\sum_{m>l=1}^M \left\| \hat{\mathbf{y}}_{t+h}^{i,(m)} - \hat{\mathbf{y}}_{t+h}^{i,(l)} \right\|_2}{M(M-1)},$$

with \mathbf{y}_{t+h}^i the target vector (wind velocity) and $\{\hat{\mathbf{y}}_{t+h}^{i,(m)}\}_{m=1}^M$ samples from the model’s predicted distribution. During training, $M = 16$ samples are considered, while $M = 100$ for testing. Training is performed for 200 epochs using Adam optimizer with an initial learning rate of 0.001 decreased at epochs 40, 60 and 120 by a factor of 0.3. Each epoch consists of 300 batches of size 32. Early

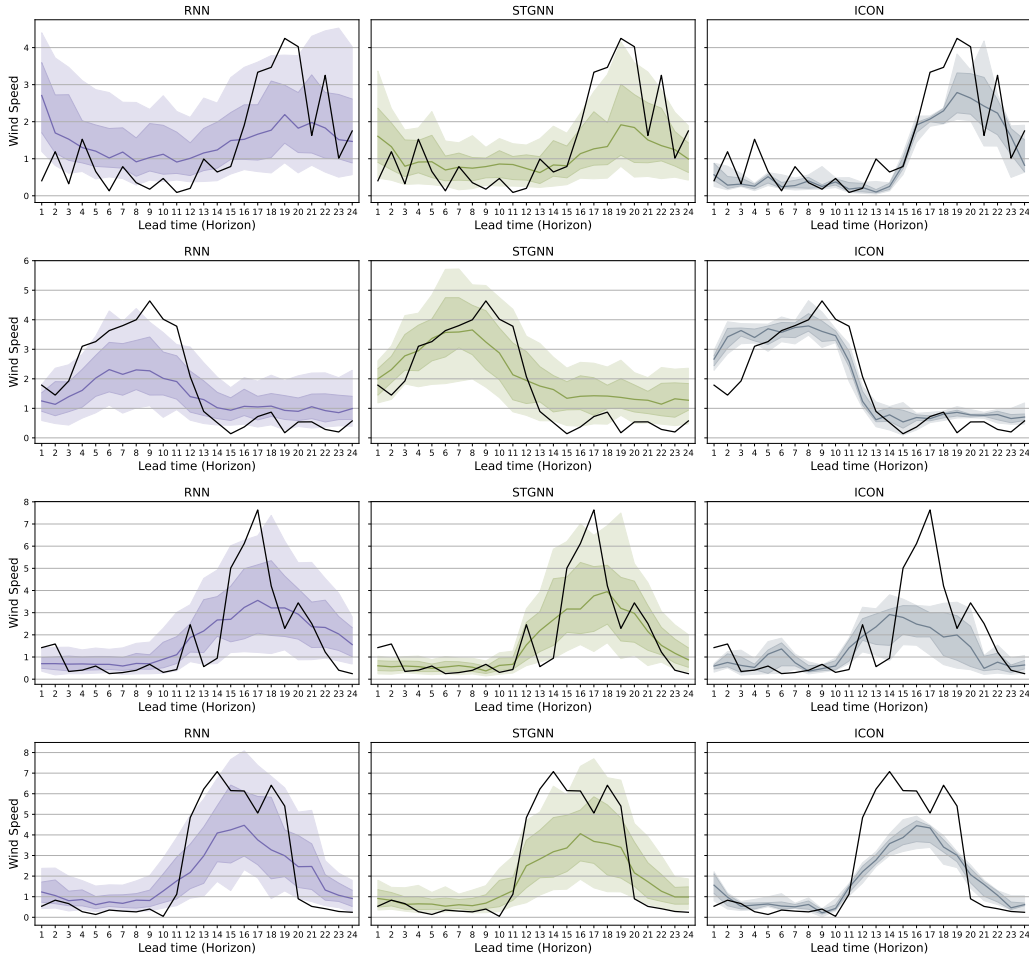
stopping is triggered after 40 epochs without improvement on the validation forecasting MAE. Models are evaluated by averaging forecasting performance over the 24h horizon and evaluating at specific lead times. To assess the quality of wind direction forecasts, wind speeds below $1m/s$ are discarded.

Hardware and Software Experiments are run on a machine with Intel(R) Xeon(R) CPU, 192 GB RAM, and NVIDIA L4 GPUs. The system ran Debian GNU/Linux 11. The code to run the experiments is developed in Python [31], by relying mainly on the following open-source libraries: PyTorch [24], PyTorch Lightning [15], PyTorch Geometric (PyG) [16], Torch Spatiotemporal [11], and the developed PeakWeather library [5] to interface with the PeakWeather dataset [4].

The source code with all the configuration files necessary to reproduce the results is open-sourced on GitHub.³ Model training and evaluation do not require particular hardware. Given the above machine and experimental settings, single model training and evaluation should complete within 3 hours.

E Illustrative forecast samples

Figure 5 displays wind speed forecasts for 24h derived from the eastward and northward wind components for some sampled forecasting windows in the test set at different station locations. The figure compares the RNN-noemb and STGNN model predictions against those from ICON-CH1-EPS. Forecast uncertainty is also visualized.



³<https://github.com/Graph-Machine-Learning-Group/peakweather-wind-forecasting>

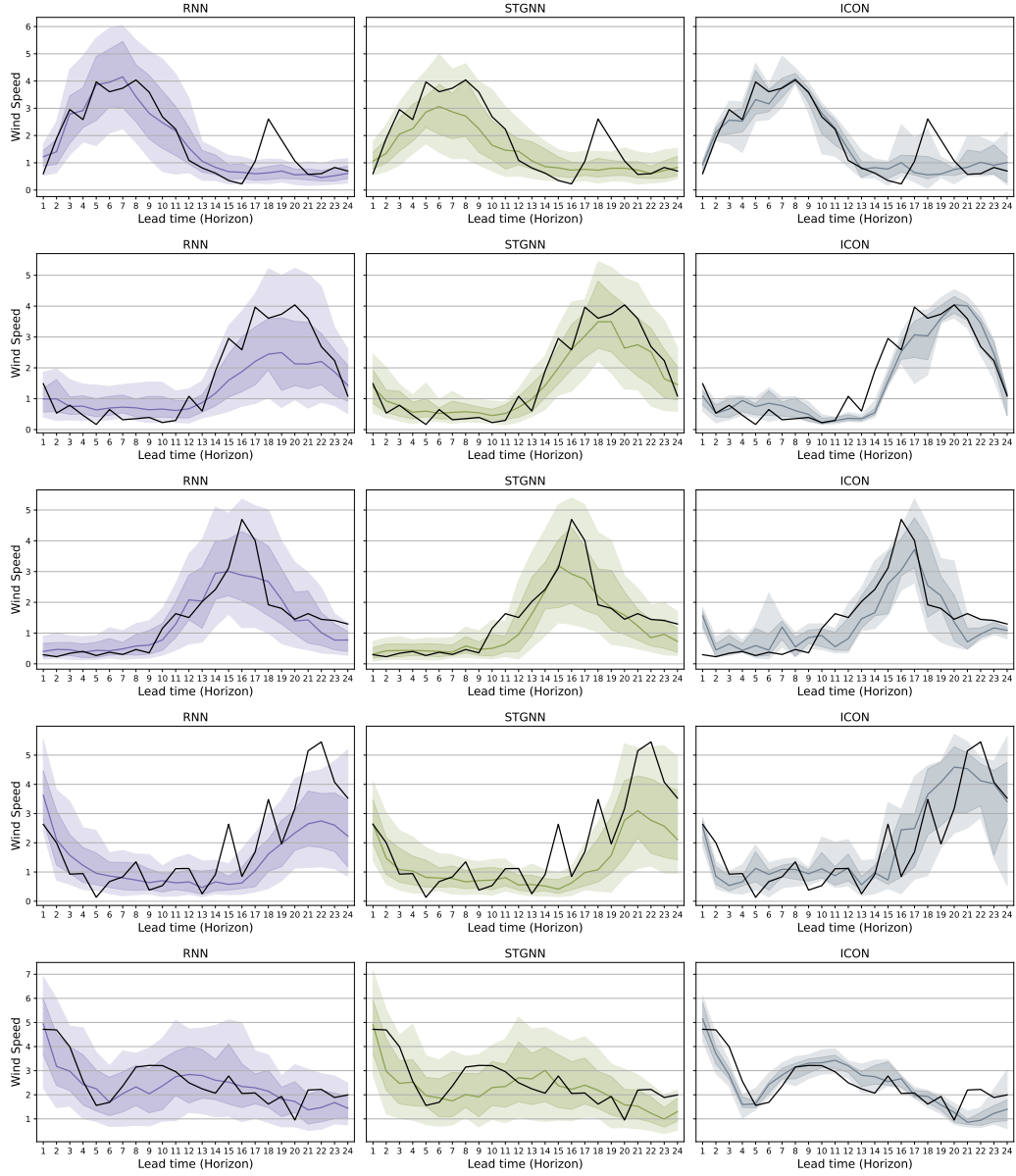


Figure 5: Comparison between RNN-noemb (left column), STGNN (center) and ICON-CH1-EPS (right column) for a selection of windows from the test set extracting the wind speed from the predicted eastward and northward wind components. For RNN-noemb and STGNN models, 50 predictive samples were produced, whereas ICON-CH1-EPS uses its 11 ensemble members. The colored solid line represents the sample median, while the black line represents the target quantity, obtained from the actual station measurements. Shaded areas denote forecast uncertainty, with the lighter region spanning from the 10th to the 90th percentile of the predictions, and the darker area highlighting the 25th to 75th percentile range.

References

- [1] DHM25, The digital height model of Switzerland, Federal Office of Topography swisstopo. <https://www.swisstopo.admin.ch/en/height-model-dhm25>. Accessed: 2025-05-10.
- [2] ECMWF Reanalysis v5 (ERA5). <https://www.ecmwf.int/en/forecasts/dataset/ecmwf-reanalysis-v5>. Accessed: 2025-05-10.
- [3] Open Data documentation of the Swiss Federal Office of Meteorology and Climatology MeteoSwiss. <https://opendatadocs.meteoswiss.ch/>. Accessed: 2025-05-10.
- [4] PeakWeather dataset. <https://huggingface.co/datasets/MeteoSwiss/PeakWeather>. Accessed: 2025-05-10.
- [5] PeakWeather Python library. <https://github.com/MeteoSwiss/PeakWeather>. Accessed: 2025-05-10.
- [6] SwissMetNet, the automatic measurement network of the Swiss Federal Office of Meteorology and Climatology MeteoSwiss. <https://www.meteoswiss.admin.ch/weather/measurement-systems/land-based-stations/automatic-measurement-network.html>. Accessed: 2025-05-10.
- [7] Peter Bauer, Alan Thorpe, and Gilbert Brunet. The quiet revolution of numerical weather prediction. *Nature*, 525(7567):47–55, 2015.
- [8] Kaifeng Bi, Lingxi Xie, Hengheng Zhang, Xin Chen, Xiaotao Gu, and Qi Tian. Pangu-weather: A 3d high-resolution model for fast and accurate global weather forecast. *arXiv preprint arXiv:2211.02556*, 2022.
- [9] David Bolton. The computation of equivalent potential temperature. *Monthly weather review*, 108(7):1046–1053, 1980.
- [10] Junyoung Chung, Caglar Gulcehre, KyungHyun Cho, and Yoshua Bengio. Empirical evaluation of gated recurrent neural networks on sequence modeling. *arXiv preprint arXiv:1412.3555*, 2014.
- [11] Andrea Cini and Ivan Marisca. Torch Spatiotemporal, 3 2022.
- [12] Andrea Cini, Ivan Marisca, Daniele Zambon, and Cesare Alippi. Taming Local Effects in Graph-based Spatiotemporal Forecasting. In A. Oh, T. Neumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors, *Advances in Neural Information Processing Systems*, volume 36, pages 55375–55393. Curran Associates, Inc., 2023.
- [13] Andrea Cini, Ivan Marisca, Daniele Zambon, and Cesare Alippi. Graph Deep Learning for Time Series Forecasting. *ACM Comput. Surv.*, 2025.
- [14] Jonathan Demaeyer, Jonas Bhend, Sebastian Lerch, Cristina Primo, Bert Van Schaeybroeck, Aitor Atencia, Zied Ben Bouallègue, Jieyu Chen, Markus Dabernig, Gavin Evans, et al. The euppbench postprocessing benchmark dataset v1. 0. *Earth System Science Data*, 15(6):2635–2653, 2023.
- [15] William Falcon and The PyTorch Lightning team. PyTorch Lightning, 3 2019.
- [16] Matthias Fey and Jan Eric Lenssen. Fast graph representation learning with pytorch geometric. *arXiv preprint arXiv:1903.02428*, 2019.
- [17] Tilmann Gneiting and Adrian E Raftery. Strictly proper scoring rules, prediction, and estimation. *Journal of the American statistical Association*, 102(477):359–378, 2007.
- [18] Rakshitha Godahewa, Christoph Bergmeir, Geoffrey I Webb, Rob J Hyndman, and Pablo Montero-Manso. Monash time series forecasting archive. *arXiv preprint arXiv:2105.06643*, 2021.

- [19] Ming Jin, Huan Yee Koh, Qingsong Wen, Daniele Zambon, Cesare Alippi, Geoffrey I. Webb, Irwin King, and Shirui Pan. A survey on graph neural networks for time series: Forecasting, classification, imputation, and anomaly detection. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 46(12):10466–10485, 2024.
- [20] Weixin Jin, Jonathan Weyn, Pengcheng Zhao, Siqi Xiang, Jiang Bian, Zuliang Fang, Haiyu Dong, Hongyu Sun, Kit Thambiratnam, and Qi Zhang. Weatherreal: A benchmark based on in-situ observations for evaluating weather models. *arXiv preprint arXiv:2409.09371*, 2024.
- [21] Remi Lam, Alvaro Sanchez-Gonzalez, Matthew Willson, Peter Wirsberger, Meire Fortunato, Ferran Alet, Suman Ravuri, Timo Ewalds, Zach Eaton-Rosen, Weihua Hu, et al. Learning skillful medium-range global weather forecasting. *Science*, 382(6677):1416–1421, 2023.
- [22] Yaguang Li, Rose Yu, Cyrus Shahabi, and Yan Liu. Diffusion convolutional recurrent neural network: Data-driven traffic forecasting. In *International Conference on Learning Representations*, 2018.
- [23] Anthony McNally, Christian Lessig, Peter Lean, Eulalie Boucher, Mihai Alexe, Ewan Pinnington, Matthew Chantry, Simon Lang, Chris Burrows, Marcin Chrust, et al. Data driven weather forecasts trained and initialised directly from observations. *arXiv preprint arXiv:2407.15586*, 2024.
- [24] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc., 2019.
- [25] Jaideep Pathak, Shashank Subramanian, Peter Harrington, Sanjeev Raja, Ashesh Chattopadhyay, Morteza Mardani, Thorsten Kurth, David Hall, Zongyi Li, Kamyar Azizzadenesheli, et al. Fourcastnet: A global data-driven high-resolution weather model using adaptive fourier neural operators. *arXiv preprint arXiv:2202.11214*, 2022.
- [26] Ilan Price, Alvaro Sanchez-Gonzalez, Ferran Alet, Tom R Andersson, Andrew El-Kadi, Dominic Masters, Timo Ewalds, Jacklynn Stott, Shakir Mohamed, Peter Battaglia, et al. Probabilistic weather forecasting with machine learning. *Nature*, 637(8044):84–90, 2025.
- [27] Stephan Rasp, Peter D Dueben, Sebastian Scher, Jonathan A Weyn, Soukayna Mouatadid, and Nils Thuerey. Weatherbench: a benchmark data set for data-driven weather forecasting. *Journal of Advances in Modeling Earth Systems*, 12(11):e2020MS002203, 2020.
- [28] Stephan Rasp, Stephan Hoyer, Alexander Merose, Ian Langmore, Peter Battaglia, Tyler Russell, Alvaro Sanchez-Gonzalez, Vivian Yang, Rob Carver, Shreya Agrawal, et al. Weatherbench 2: A benchmark for the next generation of data-driven global weather models. *Journal of Advances in Modeling Earth Systems*, 16(6):e2023MS004019, 2024.
- [29] Suman Ravuri, Karel Lenc, Matthew Willson, Dmitry Kangin, Remi Lam, Piotr Mirowski, Megan Fitzsimons, Maria Athanassiadou, Sheleem Kashem, Sam Madge, et al. Skillful precipitation nowcasting using deep generative models of radar. *Nature*, 597(7878):672–677, 2021.
- [30] Andreas Schlueter, Leonie Wagner, and Alexander Jakob Dautel. juaai/stationbench: v0.1.2, February 2025.
- [31] Guido Van Rossum and Fred L. Drake. *Python 3 Reference Manual*. CreateSpace, Scotts Valley, CA, 2009.
- [32] Stéphane Vannitsem, John Bjørnar Bremnes, Jonathan Demaeyer, Gavin R. Evans, Jonathan Flowerdew, Stephan Hemri, Sebastian Lerch, Nigel Roberts, Susanne Theis, Aitor Atencia, Zied Ben Bouallègue, Jonas Bhend, Markus Dabernig, Lesley De Cruz, Leila Hieta, Olivier Mestre, Lionel Moret, Iris Odak Plenković, Maurice Schmeits, Maxime Taillardat, Joris Van

- den Bergh, Bert Van Schaeybroeck, Kirien Whan, and Jussi Ylhäisi. Statistical postprocessing for weather forecasts: Review, challenges, and avenues in a big data world. *Bulletin of the American Meteorological Society*, 102(3):E681 – E699, 2021.
- [33] Yuchen Zhang, Mingsheng Long, Kaiyuan Chen, Lanxiang Xing, Ronghua Jin, Michael I Jordan, and Jianmin Wang. Skilful nowcasting of extreme precipitation with nowcastnet. *Nature*, 619(7970):526–532, 2023.
- [34] Xun Zhu, Yutong Xiong, Ming Wu, Gaozhen Nie, Bin Zhang, and Ziheng Yang. Weather2k: A multivariate spatio-temporal benchmark dataset for meteorological forecasting based on real-time observation data from ground weather stations. *arXiv preprint arXiv:2302.10493*, 2023.
- [35] Günther Zängl, Daniel Reinert, Pilar Rípodas, and Michael Baldauf. The icon (icosahedral non-hydrostatic) modelling framework of dwd and mpi-m: Description of the non-hydrostatic dynamical core. *Quarterly Journal of the Royal Meteorological Society*, 141(687):563–579, 2015.