# Deep Learning for Day Forecasts from Sparse Observations

Marcin Andrychowicz[*1], Lasse Espeholt[*1], Di Li[*1], Samier Merchant[2], Alexander Merose[2],
Fred Zyda[2], Shreya Agrawal[2], and Nal Kalchbrenner[*1]

[1]*Google DeepMind*
[2]*Google Research*
[*]*equal contribution*

June 2023

## Abstract

Deep neural networks offer an alternative paradigm for modeling weather conditions. The ability of neural models to make a prediction in less than a second once the data is available and to do so with very high temporal and spatial resolution, and the ability to learn directly from atmospheric observations, are just some of these models' unique advantages. Neural models trained using atmospheric observations, the highest fidelity and lowest latency data, have to date achieved good performance only up to twelve hours of lead time when compared with state-of-the-art probabilistic Numerical Weather Prediction models and only for the sole variable of precipitation. In this paper, we present MetNet-3 that extends significantly both the lead time range and the variables that an observation based neural model can predict well. MetNet-3 learns from both dense and sparse data sensors and makes predictions up to 24 hours ahead for precipitation, wind, temperature and dew point. MetNet-3 introduces a key densification technique that implicitly captures data assimilation and produces spatially dense forecasts in spite of the network training on extremely sparse targets. MetNet-3 has a high temporal and spatial resolution of, respectively, up to 2 minutes and 1 km as well as a low operational latency. We find that MetNet-3 is able to outperform the best single- and multi-member NWPs such as HRRR and ENS over the CONUS region for up to 24 hours ahead, setting a new performance milestone for observation based neural models. MetNet-3 is operational and its forecasts are served in Google Search in conjunction with other models.

## 1   Introduction

Physics based Numerical Weather Prediction (NWP) models currently drive the main forecasts that are available worldwide. These systems collect and process a large number of sparse and dense sources of observations of the atmosphere into an initial dense atmospheric representation via a process called data assimilation, which they then roll out into the future using physical laws approximations. The forward simulation is an expensive process that requires thousands of CPU hours just to make a single forecast for hours or days ahead. The spatial and temporal resolutions of the forecasts must be kept relatively low as they dramatically affect the computational cost of the simulation.

Weather models based on neural networks that use direct atmospheric observations for training offer an alternative modeling paradigm. Once the observations are available neural models have a prediction latency that is in the order of seconds. The forecast spatial resolution of the model has limited impact on computational cost that enable forecasts of one kilometer spatial resolution or higher and a very high temporal resolution in the order of minutes. Neural models can also learn atmospheric phenomena directly from the observations that capture them. This removes the need to explicitly describe a weather phenomenon using complex physics and makes it possible to model phenomena for which the physics is not well understood or that go beyond the usual domain of weather.

These advantageous properties make neural models a strong contender for an alternative paradigm for atmospheric modeling. However, high-resolution neural weather models have only been shown to perform well up to twelve hours of lead time and on the sole domain of precipitation [8]. Identifying, processing and
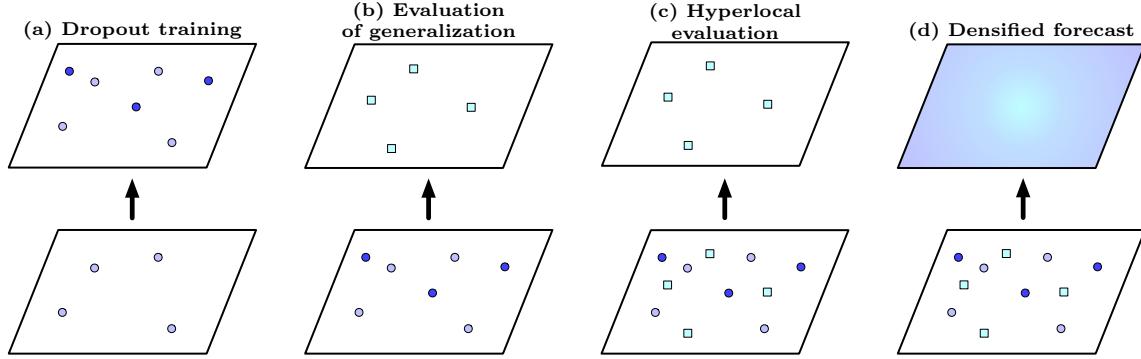
Figure 1: Abstract depiction of densification aspects. (a) During training a fraction of the weather stations are masked out from the input, while kept in the target. (b) To evaluate generalization to untrained locations, a set of weather stations represented by squares is never trained on and only used for evaluation. (c) To evaluate forecasts for the sparse locations for which data is available, these stations are fed as input during the evaluation as well. (d) The final forecasts uses the full set of training weather stations as input, and produces fully dense forecasts aided by spatial parameter sharing.

packaging for neural training the many sources of observational data that are needed to capture sufficient atmospheric information in the first place is an inordinate engineering challenge. Observational data sources come from a large number of providers with differing formats, have different spatial and temporal resolutions, and different degrees of sparsity ranging from individual points, like those from weather stations, to dense geospatial images like the observations that arise from ground-based radars and orbiting satellites. The widely different degrees of sparsity represent a novel machine learning challenge in and of themselves, as the model is expected to learn from sparse data, but produce a dense forecast.

This paper presents MetNet-3, a weather forecasting neural network that is an advance over its predecessors MetNet-1 [21] and MetNet-2 [8]. Like its predecessors, MetNet-3 maintains the same high temporal prediction frequency of 2 minutes and spatial resolution of up to 1 km. But MetNet-3 extends its lead time range from 12 hours to a full day range of 24 hours that involves dynamics well beyond extrapolation. Besides rates of precipitation, that are especially hard to predict due to their fast changing nature, MetNet-3 also predicts another set of core weather variables including surface temperature, dew point, wind speed and direction. While ground based radars provide dense precipitation measurements, observations that MetNet-3 uses for the other variables come from just 942 points that correspond to weather stations spread out across Continental United Stated (CONUS). While NWP models transform the sparse points into a dense representation during data assimilation, MetNet-3 introduces a process called *densification* to achieve this that has four main aspects (see Figure 1). The first aspect involves randomly dropping from the network's input a fraction of point observations during training, while keeping these observations as targets. The second and third aspects present two modes of evaluation of the densification, namely, evaluation on a hold-out set of stations that never appear during training to measure the network's ability to generalize spatially and perform implicit assimilation, and hyperlocal evaluation at just the specific points for which data is available. The last aspect is the inference step of densification, where the network relies on spatial parameter sharing to map all the sparse points given in the input to a fully dense image at the output, thereby producing dense forecasts.

Due to the challenge of incorporating all relevant sources of observational data that would provide a more complete picture of the recent conditions of the atmosphere, MetNet-3, like MetNet-2, still relies on an assimilated NWP initial state that describes these conditions. This state includes a dense, albeit somewhat diverging, estimate of the surface variables that MetNet-3 predicts and can aid MetNet-3 in densifying its predictions into the future for these variables.
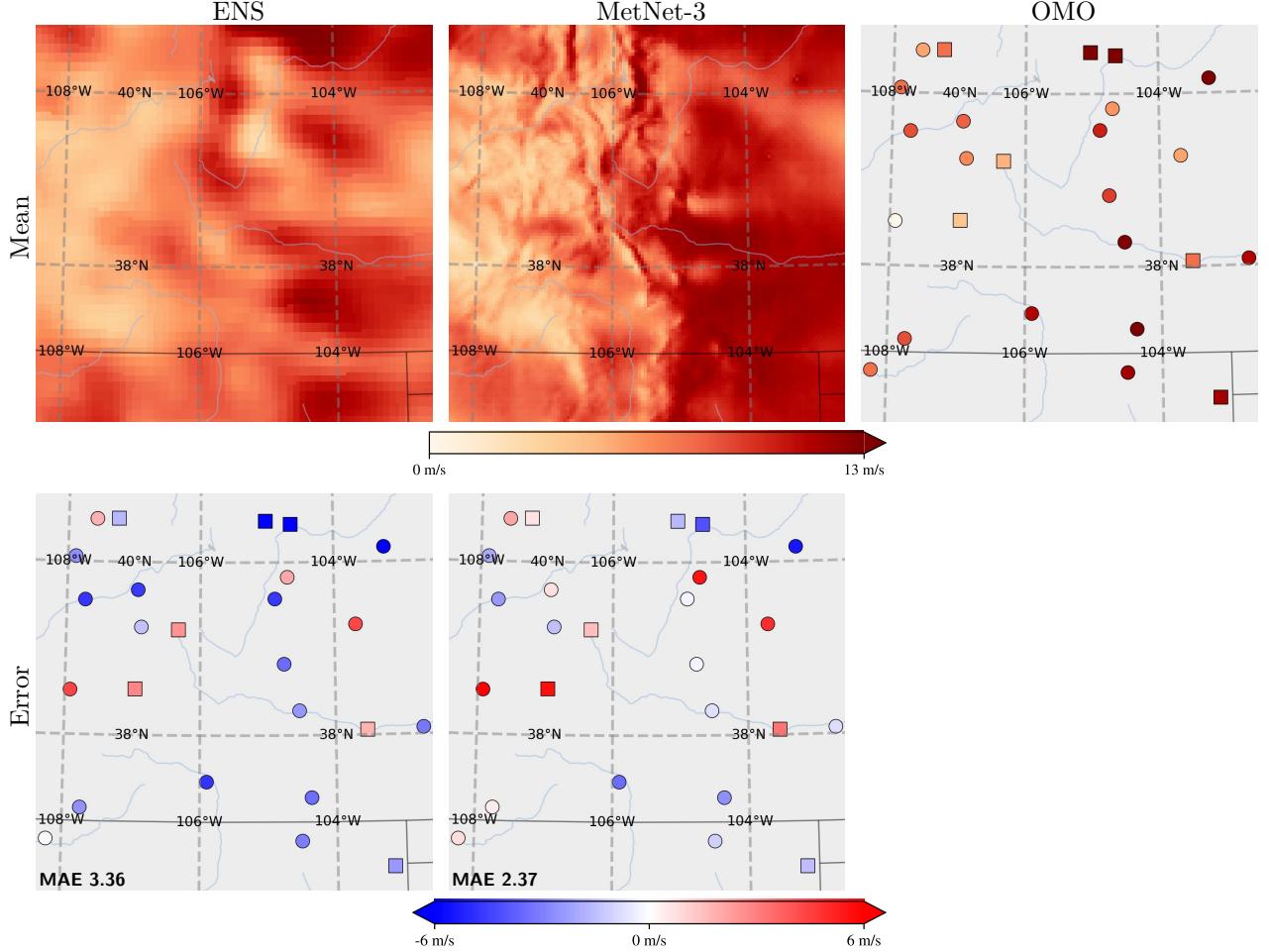
Figure 2: Case study for Sat Apr 23 2022 12:00 UTC featuring the Rocky Mountains of Colorado showing the mean of the ENS and MetNet-3 6 hour wind speed forecasts (top, left and center) along with the OMO stations ground truth (top, right) and the error of ENS and MetNet-3 on the individual weather stations (bottom). Circles and squares denote, respectively, training and test stations with MAEs calculated on both training and test stations. This example shows MetNet-3's ability to densify the targets, the higher spatial resolution of MetNet-3 as well as forecast precision on the weather stations.

## 2    Results

We evaluate MetNet-3 over CONUS on instantaneous rate of precipitation, hourly accumulated precipitation, and the surface variables: 2m temperature, 2m dewpoint, 10m wind speed and 10m wind direction.

Ground truth estimates for instantaneous precipitation come from Multi-Radar/Multi-System (MRMS) [14] and rely on radar signals. The estimates have a high temporal frequency of 2 minutes and set the base lead time frequency of MetNet-3. On the other hand, the 1-hour accumulated precipitation estimates stem from both radar signals and ground rain gauges and have a temporal frequency of 60 minutes. MRMS is generally considered a high fidelity product [24] and following [8] for evaluation we only use areas of MRMS where the radar fidelity is highest (see Supplement C).

Ground truth observations for the surface variables come from the One Minute Observations (OMO) network of weather stations [15] (see Supplement C for a map of weather stations). The weather stations include just 942 locations spread out across CONUS with observations stored for every 5th minute. MetNet-3 applies densification to this network of weather stations.

In contrast to NWPs that model uncertainty with ensemble forecasts, MetNet-3 directly outputs a
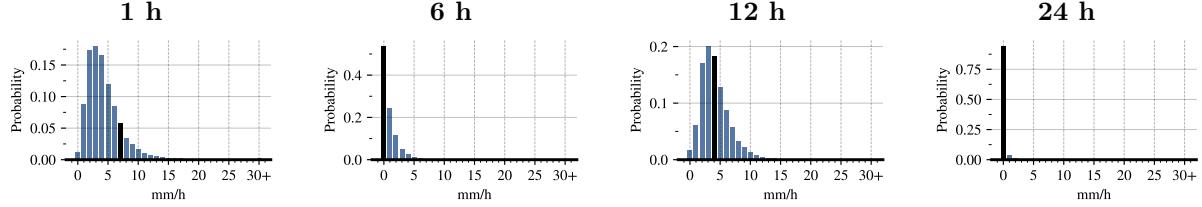
Figure 3: An example of a precipitation rate distributions from MetNet-3 forecasts for a single location for different lead times. A black colored bar indicates MRMS precipitation ground truth rate.

| Model | Region | Type | Operational Frequency | Spatial Resolution | Temporal Resolution |
|-------|--------|------|----------------------|--------------------|--------------------|
| HRRR | CONUS | Deterministic | Hourly | 3 km | Hourly |
| HREF | CONUS | Ensemble | Every 6 h | 3 km | Hourly |
| HRES | Global | Deterministic | Every 6 h | 0.1°, 7–11 km | Hourly |
| ENS | Global | Ensemble | Every 6 h | 0.2°, 14–22 km | Hourly |
| MetNet-3 | CONUS | Probabilistic | Every 10 min | 1-4 km | 2-5 min |

Figure 4: Comparison of basic characteristics of physics-based baselines used in this work and MetNet-3. MetNet-3 forecasts precipitation at 1 km / 2 min resolution and ground variables at 4 km / 5 min resolution. MetNet-3 can be run more frequently than NWP models because running the model is almost instant (about 1s for a single lead time) and requires fewer computational resources than NWP models.

marginal probability distribution for each output variable and each location using a full categorical Softmax that provides rich information beyond just the mean (see samples in Figure 3). We compare the probabilistic outputs of MetNet-3 with the outputs of advanced ensemble NWP models, including the ensemble forecast (ENS) from the European Centre for Medium-Range Weather Forecasts (ECMWF) and the High Resolution Ensemble Forecast (HREF) from the National Oceanic and Atmospheric Administration of the US (NOAA). For reference, we also include single member forecasts from the High Resolution Rapid Refresh (HRRR) and the High Resolution Forecast (HRES) by NOAA and ECMWF, respectively. We selected these models because they span the range of possible NWP models, as the former two are ensembles, while the other two are single member NWP models, and two of them are global while the other two are designed for CONUS. Figure 4 summarizes basic characteristics of the baselines used in this work and of MetNet-3. We compare the models' performance based on the metrics Continuous Ranked Probability Score (CRPS), Critical Success Index (CSI) and Mean Absolute Error (MAE). CRPS is particularly appropriate for comparison with ENS and HREF as they are ensembles of respectively 50 and 10 members and measures the accuracy of the full output distribution for all possible rates or amounts. This metric is one of the main metrics used for probabilistic forecasts and it plays an important role in guiding the development process for ENS at ECMWF [3, 7]. More details on the evaluation protocol and the metrics used can be found in Supplement D.

## 2.1 Precipitation

The first main result is that MetNet-3 obtains a higher CRPS than ENS for forecasting the rate of instantaneous precipitation over the whole lead time range of 24 hours suggesting that averaged across all rates MetNet-3's performance is superior to that of ENS (Figure 5a). When thresholding the MetNet-3 and ENS output probabilities, optimized on the categorical metric CSI, MetNet-3 outperforms ENS for the first 15 hours of lead time for light (1 mm/h) precipitation (Figure 5b) and outperforms ENS on the whole lead time range of 24 hours for heavy (8 mm/h) precipitation (Figure 5d). The skill gap between MetNet-3 and ENS is greatest in relative terms at the earliest hours and decreases gradually over time. In Figure 7, we show a 24 hour forecast of CONUS demonstrating the spread of MetNet-3 and ENS probability distributions and

**(a) CRPS**

**(b) CSI 1 mm/h**

**(c) CSI 4 mm/h**

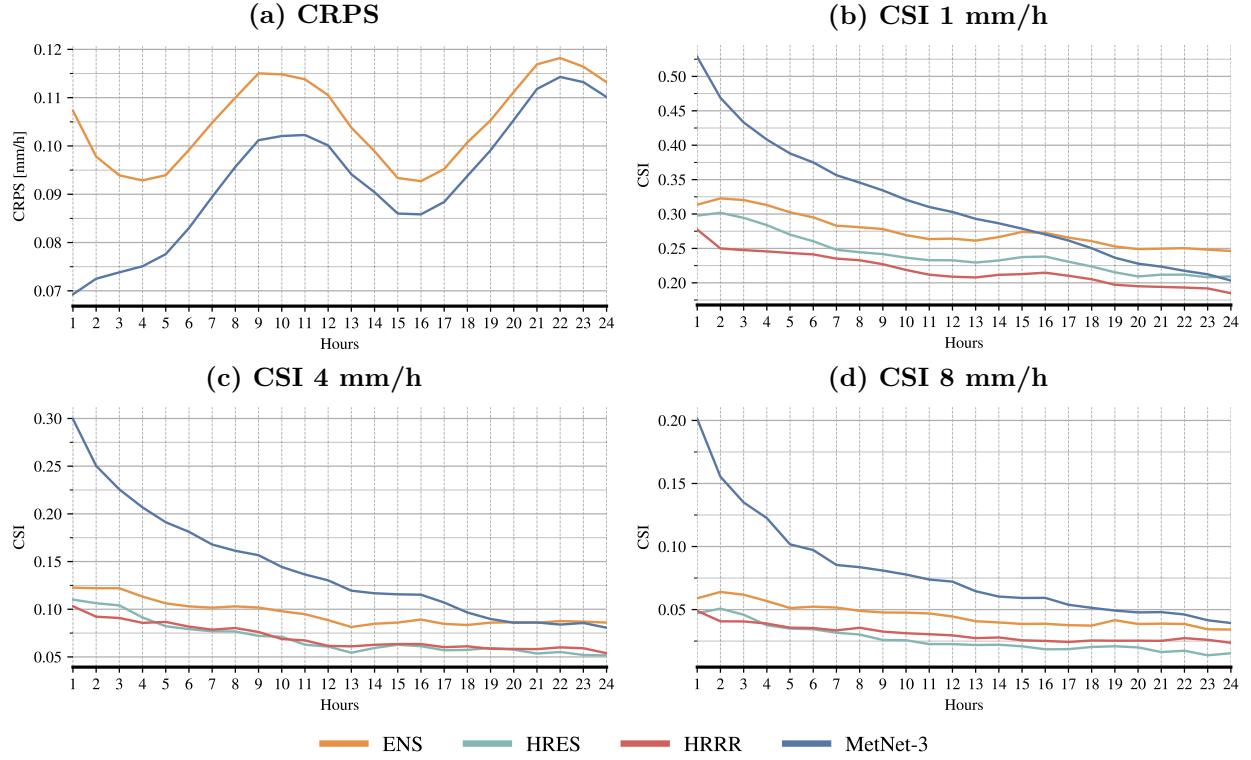**(d) CSI 8 mm/h**

ENS     HRES     HRRR     MetNet-3

Figure 5: Performance comparison between the probabilistic MetNet-3 and NWP baselines for instantaneous precipitation rate on CRPS (lower is better) and the categorical CSI (higher is better); CRPS includes all precipitation rates, whereas the CSI plots are for light (1 mm/h), moderate (4 mm/h) and heavy (8 mm/h) precipitation. Deterministic baselines (HRRR and HRES) are ommited for clarity in the CRPS plot due to performing significantly worse than the probabilistic models. Note, the thresholds for turning the probabilistic forecasts of MetNet-3 and ENS into deterministic forecasts for use in the CSI calculation, have been optimized on a validation set. HREF is omitted in all plots due to instantaneous precipitation rate being unavailable. See Supplement E for the CSI plots for other rates, the CRPS plot with the deterministic baselines included and an explanation for the CRPS lines being non-monotonic.



**(a) CRPS**

**(b) CSI 4 mm**

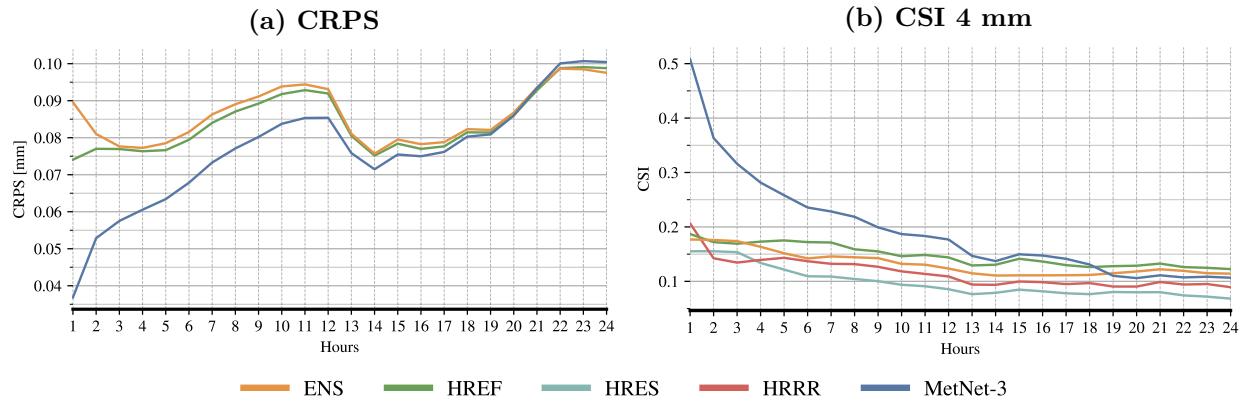ENS     HREF     HRES     HRRR     MetNet-3

Figure 6: Performance comparison between the probabilistic MetNet-3 and NWP baselines for hourly accumulated precipitation based on probabilistic CRPS (lower is better) and the categorical CSI (higher is better); CRPS includes all precipitation rates, whereas the CSI plot is for moderate (4 mm) precipitation.

the ability of MetNet-3 to predict new precipitation formations.

Hourly accumulated precipitation estimates stemming from MRMS have a frequency of only 60 minutes and provide substantially fewer distinct data frames over the same period of time that can be used for training MetNet-3 relative to those for instantaneous precipitation. Nevertheless, we find that MetNet-3 outperforms both of the multi-member NWP baselines ENS and HREF on the CRPS metric up to the first 19 hours of lead time (Figure 6a). When thresholding based on the CSI metric, MetNet-3 outperforms the baselines for the first 18 hours of lead time for moderate (4 mm/h) precipitation (Figure 6b). Detailed results can be found in the Supplement E. With these results MetNet-3 extends the lead time advantage over probabilistic NWPs for observation-based models from the 12 hours achieved by MetNet-2 all the way to 19 hours.

## 2.2    Sparse Surface Variables

Learning from weather station observations is challenging because the OMO ground truth targets are only available at 942 weather stations throughout CONUS, while weather models are required to predict weather variables at all locations. Training and evaluating on the same weather stations could lead to a situation where a model performs well on the locations that it was trained on but poorly elsewhere. To make sure that MetNet-3 performs well throughout CONUS, we apply the densification procedure for these sparse variables and we perform evaluation on a hold-out set of 20% randomly selected weather stations that are not used during MetNet-3 training and are not fed as input during evaluation either. In Figure 2, we show a case study of MetNet-3's ability to predict a densified forecast of the surface variables as well as errors on training and test weather stations.

Like for precipitation, the surface values are discretized into bins and a Softmax layer with a categorical loss is used to predict them. MetNet-3 obtains better CRPS and MAE for all surface variables than multi-member ENS for the full 24 hour range of lead times (Figure 8). In terms of CRPS that takes the full forecast distribution into account MetNet-3 shows a significant gain. For all surface variables, MetNet-3 CRPS values at all lead times up to 24 hours are better than ENS's highest CRPS, which is obtained at the shortest lead time of 1 hour ahead. This also suggests that ensemble NWPs do not model the forecast distribution particularly well in this time range. Figure 9 shows examples of MetNet-3 forecast at a single location for temperature and wind speed. Most of the observed values fall into the 80% confidence interval of MetNet-3's predicted distribution, whereas ENS's predicted distribution is very peaked and under-dispersed.

The results in terms of the MAE metric depict a similar picture where MetNet-3 achieves much better results than both the multi-member and the single-member baselines, e.g. the 20 hour MetNet-3 temperature forecast has similar MAE as a 5 hour forecast from the best performing baseline. In the hyperlocal setting, where the values of the test weather stations are given as input to the network during evaluation, the results improve further especially in the early lead times.

## 3    Discussion

MetNet-3 considers observational data as the highest fidelity target and as the ground truth for evaluation and training. An alternative would be to use assimilation or reanalysis states from NWP models as targets where the reanalysis state, as opposed to assimilation state, integrates not just past and present information about the atmosphere, but also future information [1, 4, 11, 12, 13, 16, 17, 19]. However, this option has limitations. First, we find a significant mismatch between ground truth observations and the values of the same variables provided by the NWP states. Figure 10 gives an example of the core variable of hourly accumulated precipitation in the NWP reanalysis state from ERA5 [10] and of the corresponding estimates from the gauge-corrected MRMS; the comparison shows a marked mismatch between the two. Similarly, for surface temperature, we obtain an MAE of 1.5 C and 0.9 C between OMO weather stations and the HRES and HRRR assimilation states, respectively, a margin of error that is at times larger than that of MetNet-3's forecasts themselves. See Supplement A for an example of mismatch with HRES on instantaneous precipitation. Evaluating or training against such states makes it hard to gauge the accuracy of the resulting model. A second limitation of reanalysis or assimilation states is that they are spatially and temporally coarse. ERA5 for example has a spatial resolution of approximately 25 km and a frequency of 6 hours. A strength of neural models is that their computation only grows linearly with temporal frequency and spatial resolution and using coarse targets, like ERA5, limits the potential of what neural models can learn.
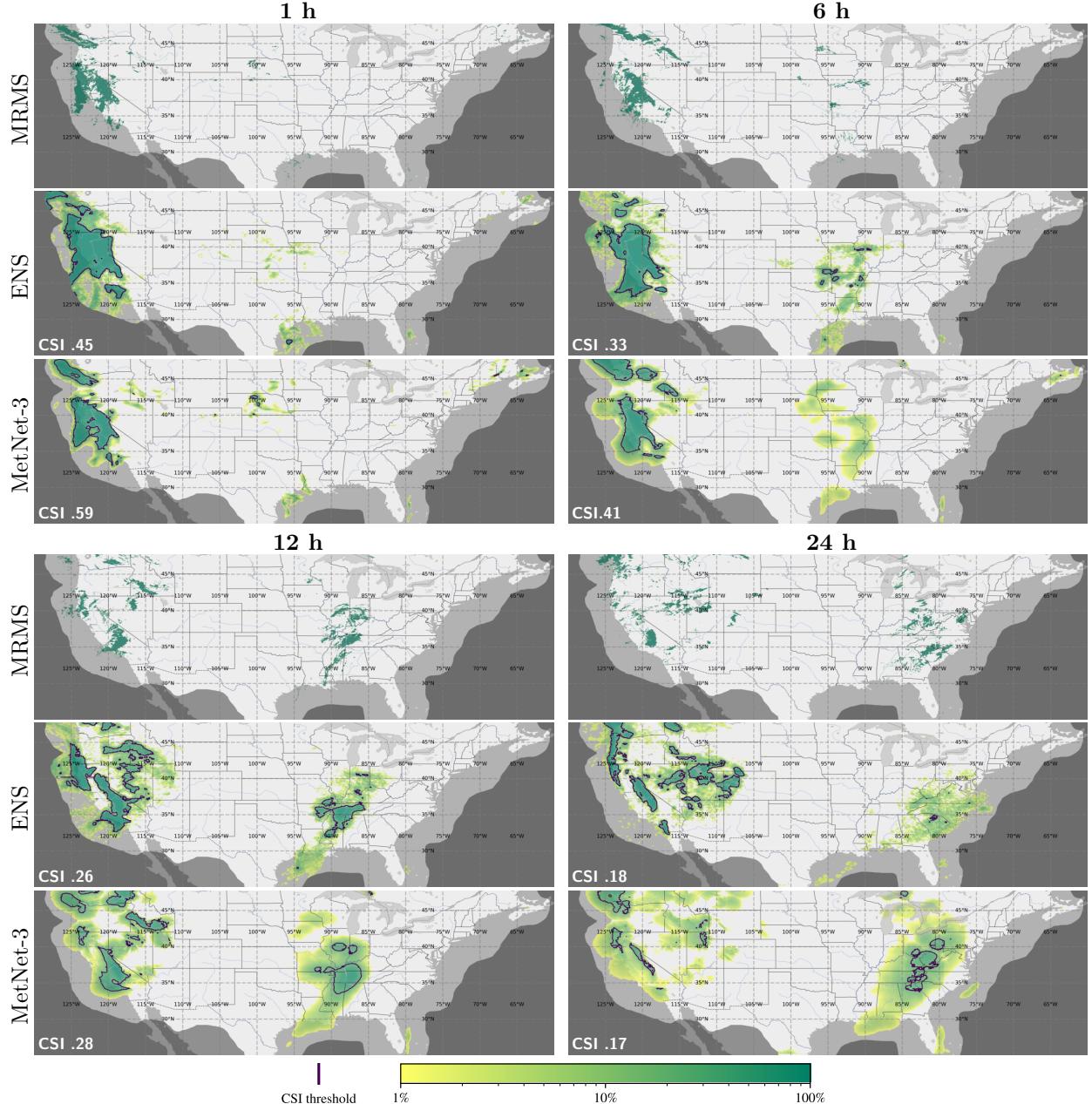
Figure 7: Case study for Thu Jan 17 2019 00:00 UTC showing the probability of instantaneous precipitation rate being above 1 mm/h on CONUS. The maps also show the prediction threshold when optimized towards CSI (dark blue contours) as well as the CSI values (lower left corners) calculated on the evaluation mask (Figure 2 in Supplement C). This specific case study shows the formation of a new large precipitation pattern in central US and *not* just extrapolation of existing patterns.
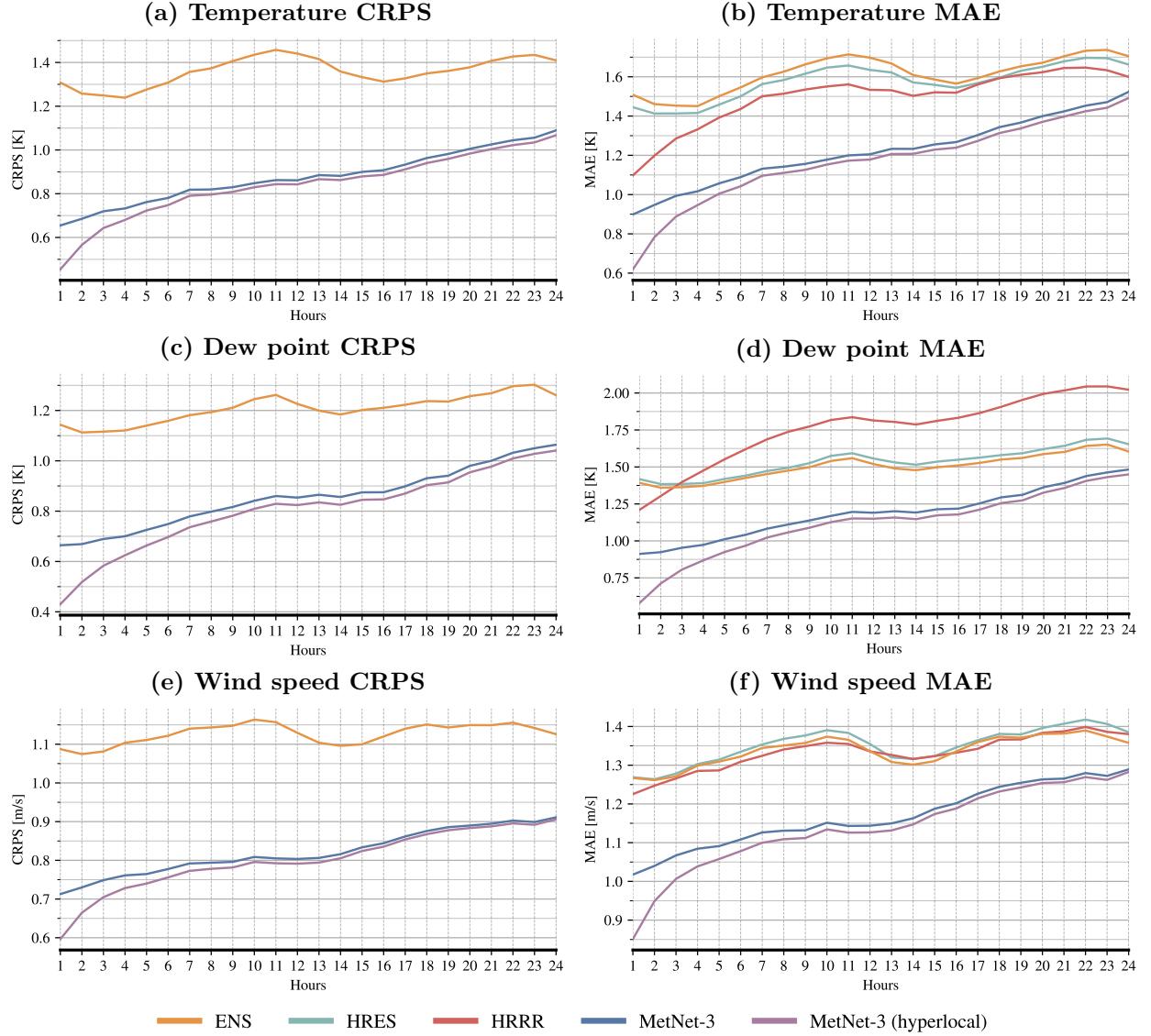
Figure 8: Performance comparison between the probabilistic MetNet-3 and NWP baselines for ground variables: temperature, dew point and wind speed based on CRPS and MAE (lower is better). Deterministic baselines (HRRR and HRES) are omitted in the CRPS plots because CRPS take the full forecast distribution into account and is therefore more appropriate for probabilistic models. Results for wind components can be found in the Supplement E. For these variables, we did not have HREF variables available.

In addition, densification from point data allows the neural model to choose an arbitrarily coarse output resolution by assigning the point data to the respective grid cell. Yet another limitation concerns the real world latency of such targets. A lag of 6 hours like that of the ENS model has a large impact on short term performance, as can be seen in Section 2. On top of that, ENS only runs 4 times per day and thus provides forecasts that rely on stale atmospheric information from up to 12 hours prior to the forecast time. This has a large effect on the operational performance of a model and for this reason MetNet-3 relies on observations with latency on the order of minutes and on the HRRR state whose latency is 55 minutes. MetNet-3 then takes another 10 minutes to generate a forecast for all of CONUS for all lead times every two minutes up to 24 hours in advance. If adjusted for operational latency, MetNet-3's gains over the NWP baselines would be larger than reported in Section 2.

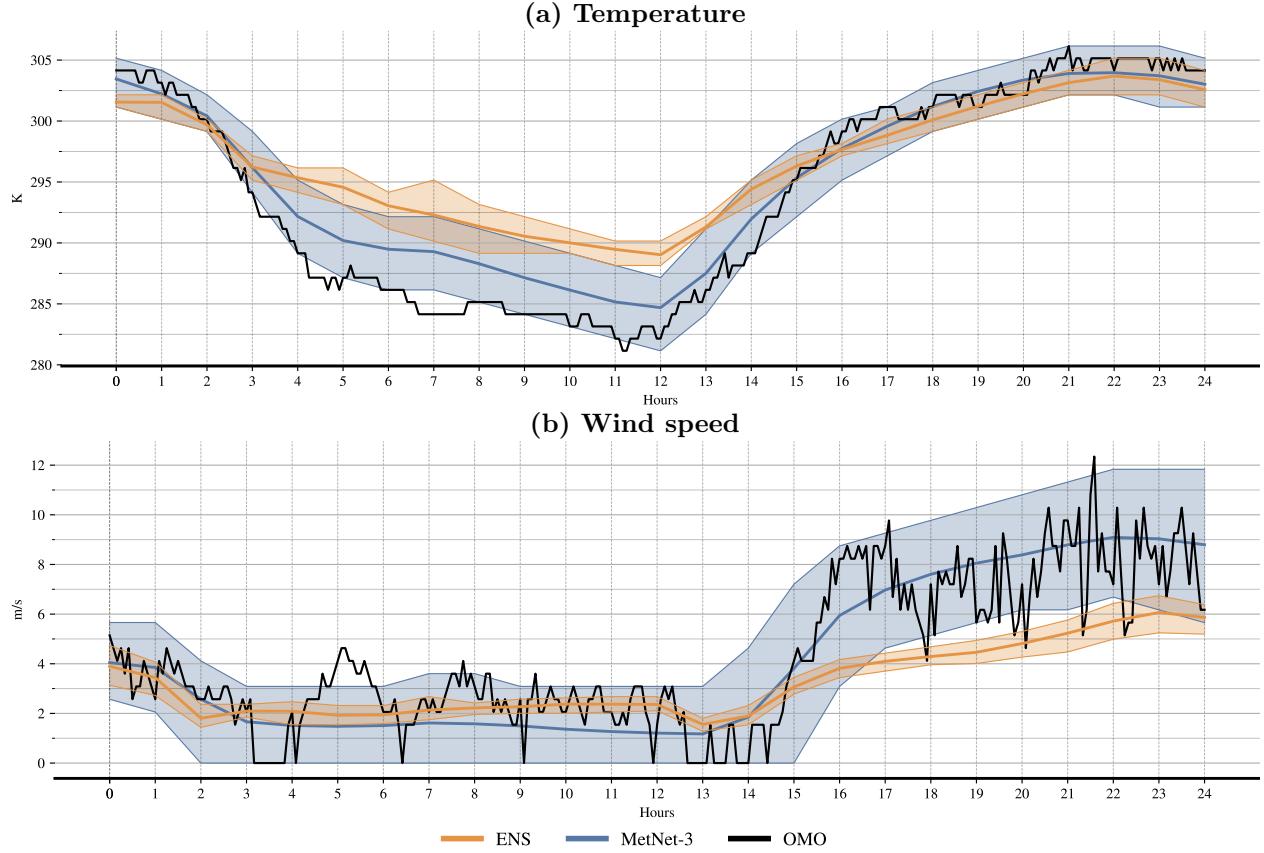When compared to MetNet-2, MetNet-3 shows a leap forward in performance. Figure 11 in Supplement E

**(a) Temperature**

**(b) Wind speed**

Figure 9: Case study for Thu Jun 10 2021 00:00 UTC comparing a MetNet-3 forecast and an ENS forecast for a single location (117.22°W, 33.91°N): Bold lines depict the means of the forecast distributions, and shaded areas correspond to 80% confidence interval based on the 10th- and 90th-quantile of the forecasted distribution in the case of MetNet-3 and the ensemble distribution in the case of ENS.

shows how the multiple innovations of MetNet-3 lead to a substantial gain. MetNet-2 in turn obtained a similar improvement over the original MetNet. This paints a picture where neural weather models keep on improving due to better architectures and observation sources. MetNet-3 still uses a tiny fraction of all available atmospheric data.

# 4 Methods

## 4.1 Dataset Creation

The data for MetNet-3 comes in input-output pairs where the inputs include radar (estimated precipitation rate and type) data from the last 90 mins, sparse OMO weather station reports from the last 6 hours, images from GOES satellites, assimilated weather state, latitude and longitude information, altitude information and current time, and outputs correspond to the future radar precipitation estimates (instantaneous radar-only precipitation rates as well as gauge-corrected hourly accumulations), measurements from ground weather stations (temperature, dew point, pressure and wind speed and direction) and assimilated weather state (the latter is only used to improve the model training and we do not treat it as ground truth). See Table 1 for more information on the inputs used, and Table 2 for more information on the targets used. The available data spans a period from July 2017 to September 2022. The training, validation and test data sets are generated without overlap from periods in sequence. Successive periods of 19 days training data, 1 day blackout, 2 day validation data, 1 day blackout, 2.5 days test data and 1 day blackout are used to sample,
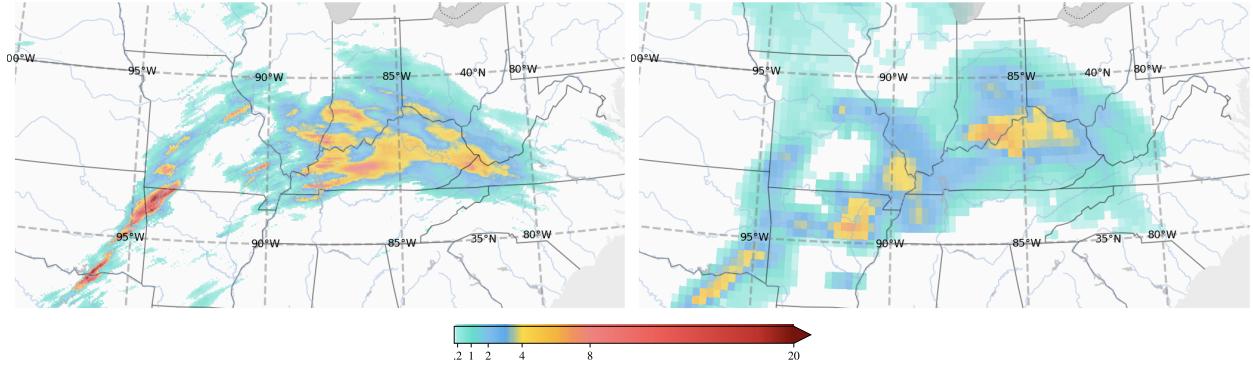
Figure 10: Hourly accumulated precipitation in millimeter accordingly to gauge-corrected MRMS product (Left) and ERA5 reanalysis data (Right), for the timestamp Sat Nov 30 2019 12:00 UTC.

respectively, training, validation and test data with no sampling in the blackout periods. To increase the number of training samples, we temporarily interpolate targets in the train split using linear interpolation whenever the observation for the exact lead time is not available. Spatially, the target patches are sampled randomly from intersections on a grid over the CONUS region spaced at .5 degrees in longitude and latitude.

For surface variables, we take the OMO station point measurements and map them to a 4 km by 4 km pixel in which the station lies. If there are multiple stations in a given region, we take the average of their measurements. For all 942 weather stations, only 12 pairs of stations are within a distance for which it is necessary to average the weather station variables. Apart from temporal splits, we also divide OMO stations into two groups: 757 training stations and 185 test stations (Supplement C, Figure 3). The data from test stations is not used in any way during training and we only report the results on the test stations. Moreover, we normally do not include past observations from the tests stations in MetNet-3 inputs even during evaluation, so that the model does not have any information about test station exact locations and produces ground forecasts representative of the full 4 km by 4 km output squares. Including past observations from the test stations allows MetNet-3 to bespoke the forecast to a particular weather station and results in hyperlocal forecasts.

## 4.2 Model and Architecture

On a high level, MetNet-3 neural network consists of three parts: topographical embeddings, U-Net [20] backbone and a MaxVit [22] transformer for capturing long-range interactions. The whole network has 227M trainable parameters.

### 4.2.1 Topographical Embeddings

It is common to feed neural weather models multiple time-independent variables containing topographical information like sea-land mask [8]. Instead of manually selecting and preparing this kind of information, we use a novel technique of *topographical embeddings*, which allows the network to automatically discover relevant topographical information and store it in embedding. More precisely, we allocate a grid of embeddings with a stride of 4 km where each point is associated with 20 scalar parameters. For each input example, we calculate the topographical embedding of each input pixel center by bilinearly interpolating the embeddings from the grid. The embedding parameters are trained together with other model parameters similarly to embeddings used in NLP.

### 4.2.2 Network Architecture

The network architecture is presented in Figure 11. The network uses two types of inputs: high-resolution, small-context (2496 km by 2496 km at 4 km resolution) ones and low-resolution, large-context ones (4992 km by 4992 km at 8 km resolution). All time slices from different high-resolution inputs (see Table 1) are

| Input | Context size | Resolution | #Channels | #Time Slices |
|---|---|---|---|---|
| Radar MRMS | 2496 km | 4 km | 2 | 11 |
| Weather stations OMO | 2496 km | 4 km | 14 | 9 |
| Elevation | 2496 km | 4 km | 1 | 1 |
| Geographical coordinates | 2496 km | 4 km | 2 | 1 |
| Topographical embeddings | 2496 km | 4 km | 20 | 1 |
| HRRR assimilation | 2496 km | 4 km | 617+1 | 1 |
| Low-resolution Radar MRMS | 4992 km | 8 km | 1 | 1 |
| GOES Satellites | 4992 km | 8 km | 16 | 1 |

Table 1: MetNet-3 spatial inputs. For HRRR assimilation, the model is given 617 channels from the assimilated state as well as one channel containing information about how stale the HRRR state is. Apart from the spatial inputs, MetNet-3 is also given the time when the prediction is made (month, day, hour and minute) and the lead time.

| Target | Source | Resolution | #Channels | Output Type | Loss Function |
|---|---|---|---|---|---|
| Precipitation | MRMS | 1 km / 2 min | 2 | Categorical | Cross Entropy |
| Surface Variables | OMO | 4 km / 5 min | 6 | Categorical | Cross Entropy |
| Assimilation | HRRR | 4 km / 1 h | 617 | Deterministic | Mean Squared Error |

Table 2: MetNet-3 targets. For precipitation, we use radar-only instanteneous precipitation estimates from MRMS as well as hourly precipitation accumulations which also take rain gauges into account. For surface variables, we use temperature, dew point and wind (speed, direction and 2 components) as reported from OMO.

first concatenated across the channel dimension, then current time is also concatenated across the channel dimension, which results in an 624 x 624 x 793 input image.

Data is then processed by a U-Net backbone, which starts with applying two convolutional ResNet blocks [9] and downsampling the data to 8 km resolution. We then pad the internal representation spatially with zeros to 4992 km by 4992 km square and concatenate with the low-resolution, large-context inputs. Afterward, we again apply two convolutional ResNet blocks and downsample the representation to 16 km resolution. Convolutional ResNet blocks can only handle local interactions and for longer lead times close to 24 hours, the targets may depend on the entire input. In order to facilitate that, we process the data at 16 km resolution using a modified version of MaxVit [22] network. MaxVit is a version of Vision Transformer (ViT, [6]) with attention over local neighbourhood as well as global gridded attention. We modify the MaxVit architecture by removing all MLP sub-blocks, adding skip connections (to the MaxVit output) after each MaxVit sub-block, and using normalized keys and queries in attention [5].

Afterwards, we take the central crop of size 768 km by 768 km, and gradually upsample the representation to 4 km resolution using skip connections from the downsampling path, at which point we again take a central crop, this time of size 512 km by 512 km. The network outputs a categorical distribution over 256 bins for each of 6 ground weather variables and a deterministic prediction for each of 617 assimilated weather state channels using an MLP with one hidden layer applied to the representation at 4 km resolution. For precipitation (both instantaneous rate and hourly accumulation), we upsample the representation to 1 km resolution and output for each pixel a categorical distribution over 512 bins. Low-level details regarding the network architecture, optimization and hyperparameters used can be found in Supplement B.
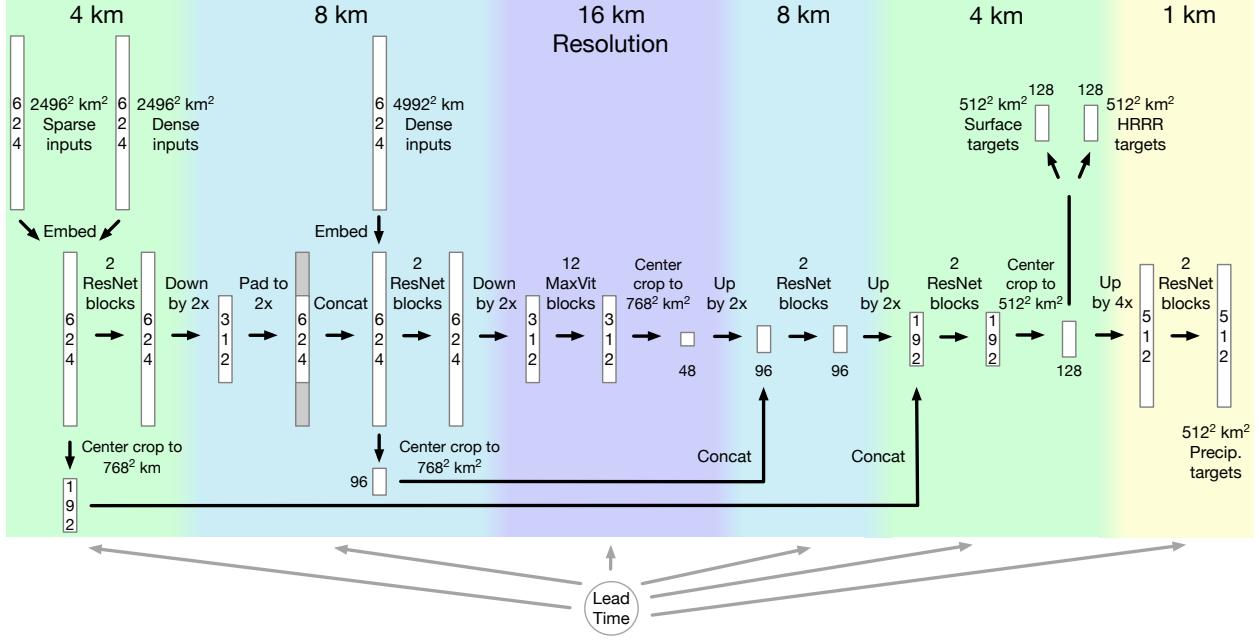
Figure 11: MetNet-3 network architecture. Rectangles denote tensors and the numbers on/under them denote their spacial sizes in pixels.

### 4.2.3 Conditioning with Lead Time

Following MetNet-2 [8], we encode the lead time as a one-hot embedding with indices from 0 to 721 representing the range between 0 and 24 hours with a 2 min interval and map them into a continuous 32-dimensional representation. Instead of feeding the lead time embedding as an input, the embedding is applied both as an additive and multiplicative factor [18] to the model inputs and to hidden representations before each activation function or self attention block. This ensures that the internal computation in the network depends directly on lead time.

The task of forecasting weather becomes significantly harder as the lead time increases which can negatively impact the model training. To counteract it, we sample the lead time during training in a biased way (exponential distribution) with $t = 24h$ being sampled 10 times less frequently than $t = 0h$. We noticed that this sampling scheme improves the results for all lead times including the long ones, which are sampled less frequently.

## 4.3 Training

The network is trained to minimize the cross-entropy loss between the ground truth data distribution and the model output. For computational efficiency, the predictions for HRRR assimilated state are deterministic and optimized with the Mean Squared Error (MSE) loss. HRRR prediction loss is included in the model solely because it improves the quality of the forecast for other variables and we do not evaluate the predictions made by the model for the assimilated state.

### 4.3.1 Densification

While we only have the ground truth for surface variables at sparse locations, the model needs to be able to generalize to all locations. To this aim, we randomly mask out each OMO station with 25% probability while training. This ensures that the model is trained to predict OMO variables even if there are no input OMO variables at the given location. (Note, this is separate from the 20% hold-out set.)

We have also noticed that there is a trade-off between the quality of precipitation and ground variables forecasts in a single model, and the results can be slightly improved by having a separate model which specializes in predicting ground variables but performs a bit worse for precipitation. Therefore, we first

train a model which is used for precipitation, and afterwards we increase the weight of the OMO loss by 100x compared to the precipitation model and finetune the model. Moreover, we disable topographical embedding (fix them to zeros) for this OMO-specific model because topographical embedding may hinder transfer between different locations, which is crucial for learning only from targets present at a sparse set of locations. See Figure 9 in Supplement E for plots comparing the two models.

### 4.3.2 Loss Scaling

As the network is trained to optimize multiple losses (cross entropy for instantaneous and accumulated precipitation rate as well as 6 OMO variables, and MSE for 617 HRRR assimilation variables) which may have very different magnitudes, it is necessary to rescale them so that their magnitudes are of similar order. Apart from using standard techniques, namely rescaling all targets for the MSE loss so that each variable has approximately mean 0 and standard deviation 1 and using manual scaling factors, we also introduce a novel technique, which relies on dynamically rescaling the gradient for each input-output sample. More precisely, after calculating the gradient of the MSE loss w.r.t. the model output for each sample, we rescale it, so that it has the same L1 norm for each output channel without changing the overall magnitude of the gradient for the sample. Let $g_{ijc}$ denote the spacial location $i, j$ and channel $c$ of the gradient w.r.t model output, and $C$ denote the number of channels. We then use the following rescaled gradient instead of $g$:

$$\hat{g}_{ijc} = \frac{C \cdot w_c}{\sum_{c'} w_{c'}} g_{ijc} \qquad\qquad w_c = \frac{1}{\sum_{i'j'} |g_{i'j'c}|} \qquad (1)$$

where the sums are over all channels ($c'$) and all spacial locations ($i', j'$) of the model output for a single input-output sample.This scaling guarantees, that the influence of each output channel is bounded and therefore even if a small fraction of the target channels are corrupted, their effect on the model is limited.

### 4.3.3 Hardware Configuration

Due to large size of the input context and internal network representations (2496 km by 2496 km at 4 km resolution and 4996 km by 4996 km at 8 km resolution), the network does not fit on a single TPU core. Instead of reducing the resolution, which could negatively impact the forecast quality, we use model parallelism. We follow MetNet-2 [8] and split the inputs, internal representation and targets into a four by four grid processed by 16 interconnected TPU cores, with each TPU core responsible for 1/16 of the area. The only exception to this rule is gridded attention in MaxVit, where we partition the data across TPU cores so that full attention windows are processed on a single core. The necessary communication at each layer is handled automatically and efficiently [2, 23].

The network is trained on 512 TPUv3 cores, where each of the 32 groups of 16 TPU cores process 2 input-output samples and the gradients from each group are synchronously aggregated after processing each batch. The fully trained MetNet-3 model took 7 days to train.

# Acknowledgements

# References

[1] 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.

[2] James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and Qiao Zhang. JAX: composable transformations of Python+NumPy programs, 2018.

[3] R. Buizza, Magdalena Alonso-Balmaseda, Andrew Brown, S.J. English, Richard Forbes, Alan Geer, T. Haiden, Martin Leutbecher, L. Magnusson, Mark Rodwell, M. Sleigh, Tim Stockdale, Frédéric Vitart, and N. Wedi. The development and evaluation process followed at ecmwf to upgrade the integrated forecasting system (ifs), 10 2018.

[4] Kang Chen, Tao Han, Junchao Gong, Lei Bai, Fenghua Ling, Jing-Jia Luo, Xi Chen, Leiming Ma, Tianning Zhang, Rui Su, et al. Fengwu: Pushing the skillful global medium-range weather forecast beyond 10 days lead. *arXiv preprint arXiv:2304.02948*, 2023.

[5] Mostafa Dehghani, Josip Djolonga, Basil Mustafa, Piotr Padlewski, Jonathan Heek, Justin Gilmer, Andreas Steiner, Mathilde Caron, Robert Geirhos, Ibrahim Alabdulmohsin, et al. Scaling vision transformers to 22 billion parameters. *arXiv preprint arXiv:2302.05442*, 2023.

[6] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. *ICLR*, 2021.

[7] ECMWF. A new tool to understand changes in ensemble forecast skill. `https://www.ecmwf.int/en/newsletter/166/news/new-tool-understand-changes-ensemble-forecast-skill`, 2021. Accessed: 2023-05-24.

[8] Lasse Espeholt, Shreya Agrawal, Casper Sønderby, Manoj Kumar, Jonathan Heek, Carla Bromberg, Cenk Gazen, Rob Carver, Marcin Andrychowicz, Jason Hickey, Aaron Bell, and Nal Kalchbrenner. Deep learning for twelve hour precipitation forecasts. *Nature Communications*, 13(1):5145, Sep 2022.

[9] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.

[10] Hans Hersbach, Bill Bell, Paul Berrisford, Shoji Hirahara, András Horányi, Joaquín Muñoz-Sabater, Julien Nicolas, Carole Peubey, Raluca Radu, Dinand Schepers, et al. The era5 global reanalysis. *Quarterly Journal of the Royal Meteorological Society*, 146(730):1999–2049, 2020.

[11] Ryan Keisler. Forecasting global weather with graph neural networks. *arXiv preprint arXiv:2202.07575*, 2022.

[12] Thorsten Kurth, Shashank Subramanian, Peter Harrington, Jaideep Pathak, Morteza Mardani, David Hall, Andrea Miele, Karthik Kashinath, and Animashree Anandkumar. Fourcastnet: Accelerating global high-resolution weather forecasting using adaptive fourier neural operators. *arXiv preprint arXiv:2208.05419*, 2022.

[13] Remi Lam, Alvaro Sanchez-Gonzalez, Matthew Willson, Peter Wirnsberger, Meire Fortunato, Alexander Pritzel, Suman Ravuri, Timo Ewalds, Ferran Alet, Zach Eaton-Rosen, et al. Graphcast: Learning skillful medium-range global weather forecasting. *arXiv preprint arXiv:2212.12794*, 2022.

[14] MRMS. Multi-radar/multi-sensor system (mrms). `https://www.nssl.noaa.gov/projects/mrms/`, 2021. Accessed: 2021-06-01.

[15] NCEP. Ncep - 1-minute asos data. `https://madis.ncep.noaa.gov/madis_OMO.shtml`, 2017. Accessed: 2023-05-16.

[16] Tung Nguyen, Johannes Brandstetter, Ashish Kapoor, Jayesh K Gupta, and Aditya Grover. Climax: A foundation model for weather and climate. *arXiv preprint arXiv:2301.10343*, 2023.

[17] 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.

[18] Ethan Perez, Florian Strub, Harm De Vries, Vincent Dumoulin, and Aaron Courville. Film: Visual reasoning with a general conditioning layer. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, 2018.

[19] 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.

[20] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation, 2015. cite arxiv:1505.04597Comment: conditionally accepted at MICCAI 2015.

[21] Casper Kaae Sønderby, Lasse Espeholt, Jonathan Heek, Mostafa Dehghani, Avital Oliver, Tim Salimans, Shreya Agrawal, Jason Hickey, and Nal Kalchbrenner. Metnet: A neural weather model for precipitation forecasting. *arXiv preprint arXiv:2003.12140*, 2020.

[22] Zhengzhong Tu, Hossein Talebi, Han Zhang, Feng Yang, Peyman Milanfar, Alan Bovik, and Yinxiao Li. Maxvit: Multi-axis vision transformer. *ECCV*, 2022.

[23] Yuanzhong Xu, HyoukJoong Lee, Dehao Chen, Blake A. Hechtman, Yanping Huang, Rahul Joshi, Maxim Krikun, Dmitry Lepikhin, Andy Ly, Marcello Maggioni, Ruoming Pang, Noam Shazeer, Shibo Wang, Tao Wang, Yonghui Wu, and Zhifeng Chen. GSPMD: general and scalable parallelization for ML computation graphs. *CoRR*, abs/2105.04663, 2021.

[24] Jian Zhang, Kenneth Howard, Carrie Langston, Brian Kaney, Youcun Qi, Lin Tang, Heather Grams, Yadong Wang, Stephen Cocks, Steven Martinaitis, Ami Arthur, Karen Cooper, Jeff Brogden, and David Kitzmiller. Multi-radar multi-sensor (mrms) quantitative precipitation estimation: Initial operating capabilities. *Bulletin of the American Meteorological Society*, 97(4):621 – 638, 2016.

# Supplemental Material to
# Deep Learning for Day Forecasts
# from Sparse Observations

Marcin Andrychowicz[*1], Lasse Espeholt[*1], Di Li[*1], Samier Merchant[2], Alexander Merose[2],
Fred Zyda[2], Shreya Agrawal[2], and Nal Kalchbrenner[*1]

[1]*Google DeepMind*
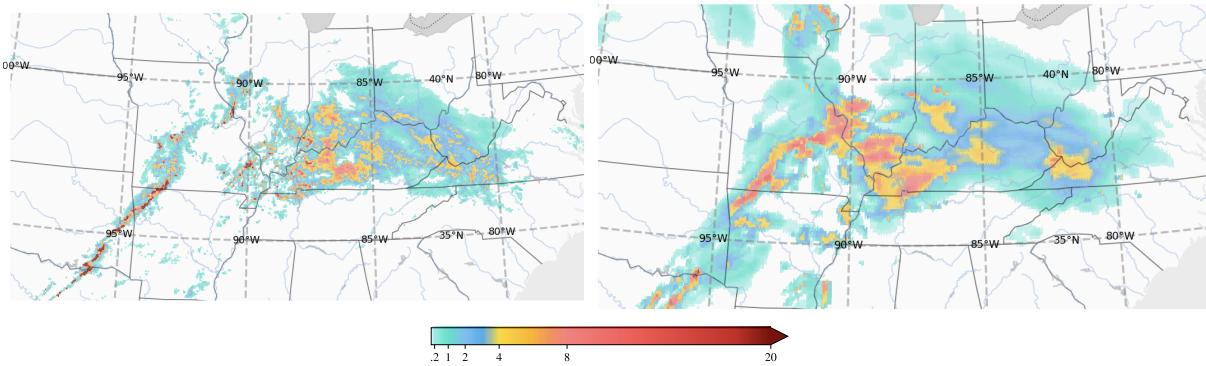[2]*Google Research*
[*]*equal contribution*

June 2023

## A   Data



Figure 1: Precipitation rate in mm/h accordingly to MRMS (Left) and HRES assimilation (Right), for the timestamp Sat Nov 30 2019 12:00 UTC.

## B   Supplement: Model and Training

Optimization hyperparameters can be found in Table 1. Below we list additional technical details related to the network architecture:

**Inputs**   The high-resolution MRMS input has two channels — instantaneous precipitation rate and precipitation type, while the low-resolution MRMS input only contains the precipitation rate. Precipitation rate inputs are preprocessed using the following transformation: $\tanh(\log(r+1)/4)$, where $r$ is the precipitation rate in mm/h. All other input channels are normalized to have mean and standard deviation values that are approximately 0 and 1, respectively. We use time slices with the following offsets (in minutes) — high-resolution MRMS: -90, -75, -60, -45, -30, -25, -20, -15, -10, -5, 0; OMO: -360, -180, -120, -60, -30, -15, -10, -5, 0; all other inputs: 0. Inputs are embedded to the internal representation of size 512 using a linear layer.

1

**Network**  We use 512 channels throughout the whole network with the exception of 2 MLPs (one at 4 km resolution and one at 1 km resolution) which produce the network outputs which have a single hidden layer of size 4096. All convolutions have kernels of size (3, 3) and are not dilated. For computational efficiency, we use mixed precision [12] with most of the computation performed in bfloat16 format.

**MaxVit**  We use 12 modified MaxVit [16] blocks. We introduced the following modifications compared to the original architecture: we removed MLP sub-blocks which were present in the original MaxVit architecture, we use normalized keys and queries [3] and we introduce skip connections from the output of each MaxVit block to the output of MaxVit. More precisely, the final output of MaxVit is a linear transformation of the outputs after each sub-block (after summing with the residual branch). All attention windows have size 8 by 8 and we use 32 attention heads. MBConv [7] in MaxVit uses the expansion rate of 4, and squeeze-and-excitation (SE, [8]) with the bottleneck ratio of 0.25.

**U-Net**  In the downsampling path of U-Net, we apply 2 convolutional ResNet blocks and downsample by 2x with max pooling on 4 km and 8 km resolution levels. In the upsampling path, we upsample using a transposed convolution [10] with kernel (2, 2) and stride (2, 2) on both 16 km and 8 km level, and then apply 2 convolutional ResNet blocks. Upsampling from 4 km to 1 km resolution is performed by repeating each activation across a 4 by 4 pixels square and applying again 2 ResNet blocks.

**Normalization**  We use pre-activation (pre-LN, [17]) layer normalization [1] throughout the network. We also apply layer normalization after each convolution which is not the last convolution in the given sub-block.

**Lead Time Conditioning**  We use additive-multiplicative conditioning (FiLM, [13]) on lead time throughout the network. The conditioning is applied to the the network inputs and after each layer normalization. All additive and multiplicative factors are outputted by a single MLP with one hidden layer of size 32 which takes as input one-hot encoded lead time. The second layer of this MLP is initialized so that at initialization the conditioning is an identity function.

**Topographical Embeddings**  To limit the number of parameters in the topographical embeddings, we only allocate topographical embeddings for the region 14.8-59.9N, 150.7-39.3W. This results in 3M points on a grid with a stride of 4 km and 60M trainable parameters for embeddings of size 20.

**Activation Functions**  We use GELU [5] inside MBConv (in MaxVit) and ReLu in all other places.

**Initialization**  We use LecunNormal initializer. Additionally we rescale the initialization of the last linear layer in each sub-block in MaxVit by $1/\sqrt{N}$, where N is the number of sub-blocks those outputs are added on the given residual connection as described in [2].

**Regularization**  We apply Dropout [15] with the rate of 0.1 before adding the output of each sub-module to the residual branch and after the first convolution in each ResNet block. We use stochastic depth [9] in MaxVit with the probability of dropping a given sub-module (i.e. MBConv, local attention or gridded attention) increasing linearly thorough the network from 0 to 0.2. We also use weight decay coefficient 0.1 as defined in AdamW [11].

| Training Hyperparameters | Value |
|---|---|
| Optimizer | AdamW [11] |
| Learning rate | 8e-5 |
| AdamW $\beta_1$ | 0.9 |
| AdamW $\beta_2$ | 0.999 |
| Weight Decay | 0.1 |
| Polyak Decay | 0.9999 |
| Batch size | 64 |
| Training steps | 260k |
| OMO finetuning steps | 80k |

Table 1: Optimization hyperparameters for MetNet-3.

## B.1    Outputs, Targets and Losses

| Target | Resolution | #Channels | Loss Function | #Bins | Bin Size |
|---|---|---|---|---|---|
| MRMS rate | 1 km | 1 | Cross Entropy | 512 | 0.2 mm/h |
| MRMS accumulation | 1 km | 1 | Cross Entropy | 512 | 0.2 mm |
| OMO temperature | 4 km | 1 | Cross Entropy | 256 | 1 K |
| OMO dew point | 4 km | 1 | Cross Entropy | 256 | 1 K |
| OMO wind speed | 4 km | 1 | Cross Entropy | 256 | 0.1 knot |
| OMO wind components | 4 km | 2 | Cross Entropy | 256 | 0.1 knot |
| OMO wind direction | 4 km | 1 | Cross Entropy | 180 | 2 degrees |
| HRRR assimilation | 4 km | 617 | MSE | N/A | N/A |

Table 2: Details of outputs produced by MetNet-3.

Table 2 lists different outputs produced by MetNet-3. As the network is trained to optimize multiple losses, which may have very different magnitudes, it is necessary to rescale them so that their magnitudes are of similar order of magnitude. To this aim, we first rescale all targets for the MSE loss so that each variable has approximately mean 0 and standard deviation 1, and apply dynamic gradient rescaling described in the main article.

We also introduce additional manual scaling factors:

- HRRR loss is multipled by 10 and divided by the number of HRRR channels being predict (617).

- We additionally increase the weight on HRRR channels corresponding to OMO ground variables the model predicts, namely 2m temperature, 2m dew point and 10m wind components by 30x. The weight of the remaining channels is decreased so that this step does not change the average weight of a HRRR channel.

- Each OMO target channel has the same weight with the sum of their weights being set to 0.01 for the standard (precipitation) model and increased to 1 for the OMO model finetuning.

# C    Supplement: Evaluation

**Non-monotonic CRPS plots**    We filter our evaluation dataset to only include locations and times when historical forecasts for all baselines are available. In particular, historical ENS forecasts are only available for two runs per day (00 and 12 UTC) so all our evaluations only start at two times during the day. Because

of that, the expected amount of precipitation depends on the lead time. Higher amounts of precipitation generally result in higher CRPS values, which results in cases when CRPS counter-intuitively decreases with lead time. We do not observe a similar phenomenon on CSI plots, because CSI scores are by definition normalized (Supplement D.1).

**MRMS and HRRR mask**  The quality of MRMS data varies between locations depending mostly on the distance from the nearest radar. While we use all available data for training, we only evaluate using data from locations with the highest quality of radar data (Figure 2).



Figure 2: Training (Left) and evaluation (Right) masks used for MRMS and HRRR targets.

**OMO ground truth**  This figure represents the OMO network of weather stations, also known as 1-minute FAA Automated Surface Observing System (ASOS) or formerly high-frequency METAR.
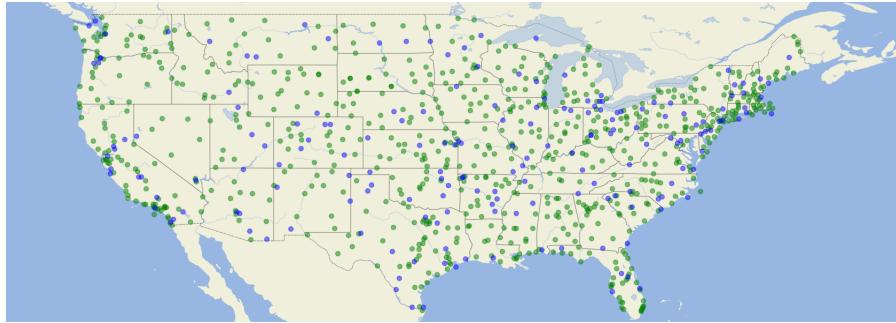


Figure 3: OMO weather stations. Blue are test stations, and green are training stations.

# D  Supplement: Evaluation Metrics

We evaluate the quality of the forecasts using three different metrics, the Continuous Ranked Probability Score (CRPS) [6], the Critical Skill Index (CSI) [14], and the Mean Absolute Error (MAE).

## D.1  Critical Success Index (CSI)

The CSI score is a binary categorical score which we use to evaluate the quality of precipitation forecasts.

$$CSI = TP/(TP + FN + FP) \tag{1}$$

where TP are true positives, FN are false negatives and FP are false positives. The CSI score is not directly applicable to the probability distributions that MetNet-3 or ensemble baselines (HREF and ENS) produce. To make a categorical decision, for a binary category corresponding to an amount of precipitation greater or equal to a given rate $r$, we calculate on a validation held-out set a probability threshold between 0 and 1 which maximizes CSI separately for each lead time. If the total predicted probability mass for rates $\geq r$ exceeds the threshold, then we take it to be a positive prediction for this rate category. This is the same procedure as used in MetNet-2 [4].

We choose CSI over similar metrics for binary classification, because it disregards the number of true negatives, i.e. the cases when there was no precipitation (or the precipitation rate was below the specified

evaluation rate) and the model predicted that correctly, and in the case of precipitation the vast majority of cases are of this type.

## D.2   Continuous Ranked Probability Score (CRPS)

CRPS in essence is the mean squared error between the cumulative density function (CDF) of the prediction and that of the ground truth integrated over the whole range of possible values. We calculate it on the discretized set of values, i.e.

$$CRPS = \sum_{i=1}^{N} (P_M(y \leq u_i) - \mathbb{1}(y \leq u_i))^2 \times \texttt{bin size}, \tag{2}$$

where $i$ iterates over all discretization bins, $u_i$ if the upper end of the $i$-th bin, $y$ is the ground truth and $P_M(y \leq u_i)$ denotes the probability that $y \leq u_i$ under the model.

## D.3   Mean Absolute Error (MAE)

Mean Absolute Error (MAE) is defined as

$$MAE = |\hat{y} - y|,$$

where $y$ if the ground truth and $\hat{y}$ is the deterministic prediction. For MetNet-3 and ensemble baselines (HREF and ENS) we take the median of the forecast distribution as $\hat{y}$. We choose median, and not mean, because median minimizes MAE for a perfect model. MAE is not a suitable metric for very skewed distributions, and therefore we do no apply it to precipitation.

# E   Supplement: Additional Results

In this section we present some additional results:

- Fig. 4: CRPS plots for precipitation including deterministic baselines (HRRR and HRES).

- Fig. 5: Instantaneous precipitation rate CSI plots for additional rates.

- Fig. 6: Hourly accumulated precipitation CSI plots for additional rates.

- Fig. 7–8: Results for surface wind U, V components.

- Fig. 9: Comparison of the standard version of MetNet-3 and the one finetuned for improved performance on ground variables.

- Fig. 10: Ablations with topographical embeddings and large-context inputs removed.

- Fig. 11: Comparison between MetNet-2 and MetNet-3.

# References

[1] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016.

[2] Rewon Child, Scott Gray, Alec Radford, and Ilya Sutskever. Generating long sequences with sparse transformers. *arXiv preprint arXiv:1904.10509*, 2019.

[3] Mostafa Dehghani, Josip Djolonga, Basil Mustafa, Piotr Padlewski, Jonathan Heek, Justin Gilmer, Andreas Steiner, Mathilde Caron, Robert Geirhos, Ibrahim Alabdulmohsin, et al. Scaling vision transformers to 22 billion parameters. *arXiv preprint arXiv:2302.05442*, 2023.
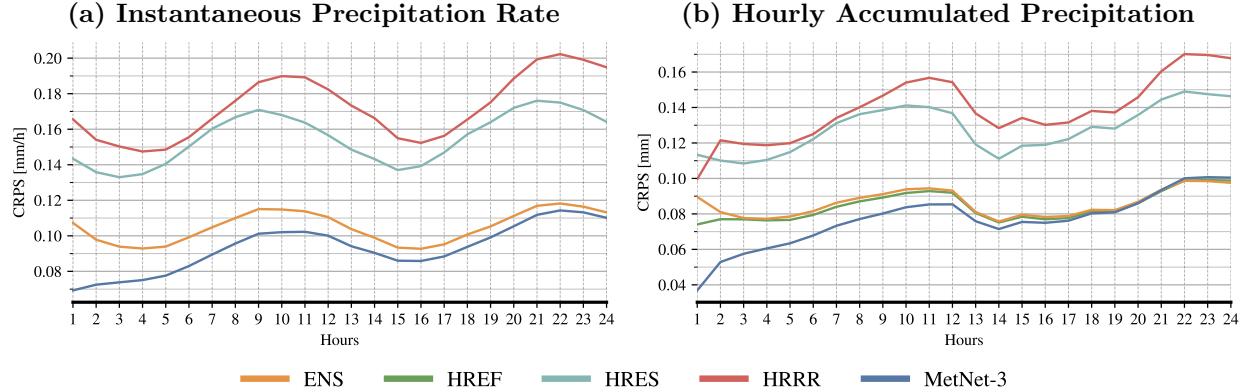
**(a) Instantaneous Precipitation Rate**    **(b) Hourly Accumulated Precipitation**



Figure 4: CRPS values for MetNet-3 and baselines based on different precipitation measurements.

[4] Lasse Espeholt, Shreya Agrawal, Casper Sønderby, Manoj Kumar, Jonathan Heek, Carla Bromberg, Cenk Gazen, Rob Carver, Marcin Andrychowicz, Jason Hickey, Aaron Bell, and Nal Kalchbrenner. Deep learning for twelve hour precipitation forecasts. *Nature Communications*, 13(1):5145, Sep 2022.

[5] Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus). *arXiv preprint arXiv:1606.08415*, 2016.

[6] Hans Hersbach. Decomposition of the continuous ranked probability score for ensemble prediction systems. *Weather and Forecasting*, 15(5):559 – 570, 2000.

[7] Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861*, 2017.

[8] Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7132–7141, 2018.

[9] Gao Huang, Yu Sun, Zhuang Liu, Daniel Sedra, and Kilian Q Weinberger. Deep networks with stochastic depth. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part IV 14*, pages 646–661. Springer, 2016.

[10] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3431–3440, 2015.

[11] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.

[12] Paulius Micikevicius, Sharan Narang, Jonah Alben, Gregory Diamos, Erich Elsen, David Garcia, Boris Ginsburg, Michael Houston, Oleksii Kuchaiev, Ganesh Venkatesh, et al. Mixed precision training. *arXiv preprint arXiv:1710.03740*, 2017.

[13] Ethan Perez, Florian Strub, Harm De Vries, Vincent Dumoulin, and Aaron Courville. Film: Visual reasoning with a general conditioning layer. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, 2018.

[14] R.J. Donaldson, R.M. Dyer, and M.J. Kraus. An objective evaluator of techniques for predicting severe weather events. In *Preprints, Ninth Conference on Severe Local Storms*, Norman, OK USA, 1975. American Meteorological Society.

[15] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *The journal of machine learning research*, 15(1):1929–1958, 2014.
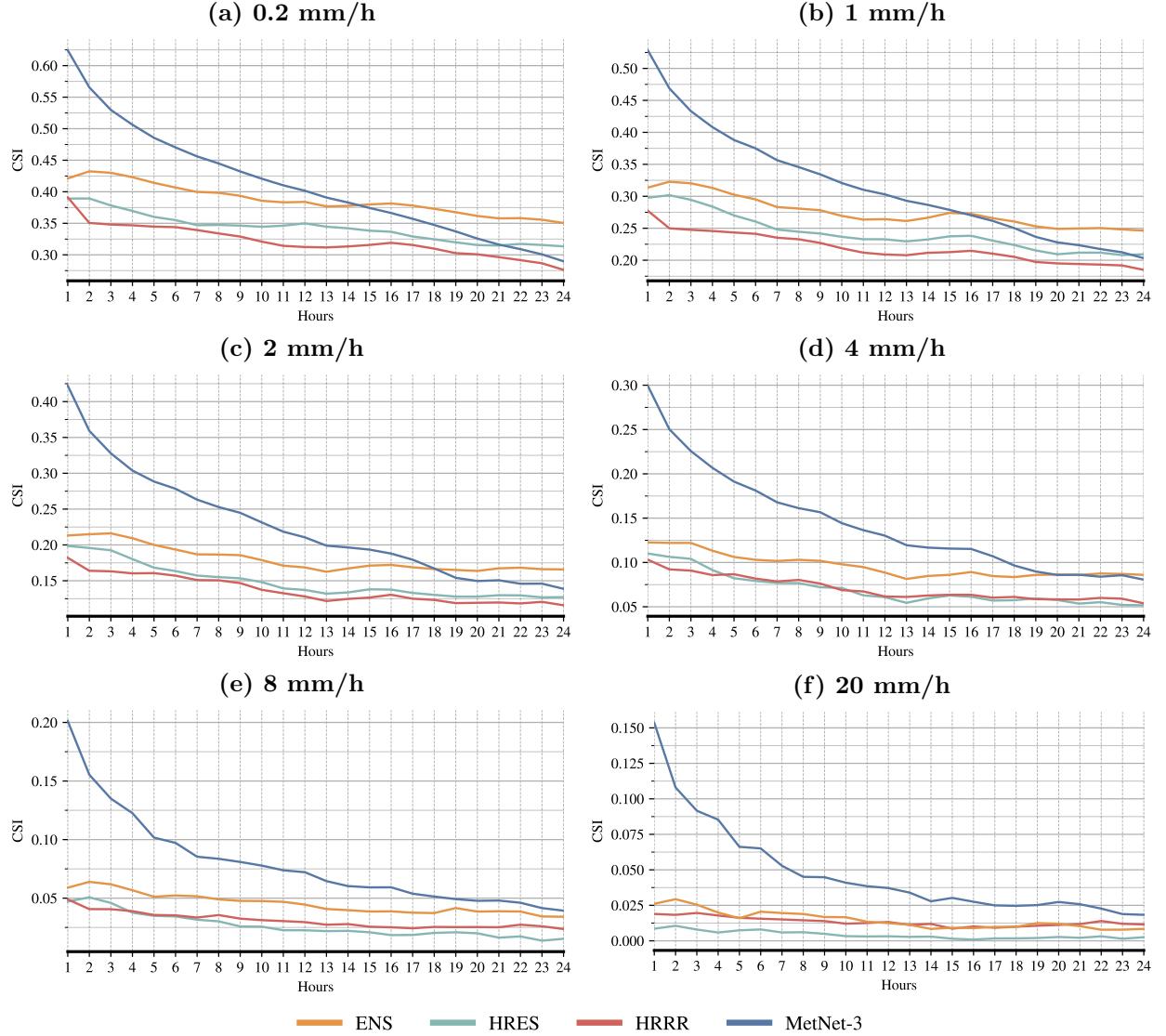
Figure 5: CSI values for instantaneous precipitation rate.

[16] Zhengzhong Tu, Hossein Talebi, Han Zhang, Feng Yang, Peyman Milanfar, Alan Bovik, and Yinxiao Li. Maxvit: Multi-axis vision transformer. *ECCV*, 2022.

[17] Ruibin Xiong, Yunchang Yang, Di He, Kai Zheng, Shuxin Zheng, Chen Xing, Huishuai Zhang, Yanyan Lan, Liwei Wang, and Tieyan Liu. On layer normalization in the transformer architecture. In *International Conference on Machine Learning*, pages 10524–10533. PMLR, 2020.
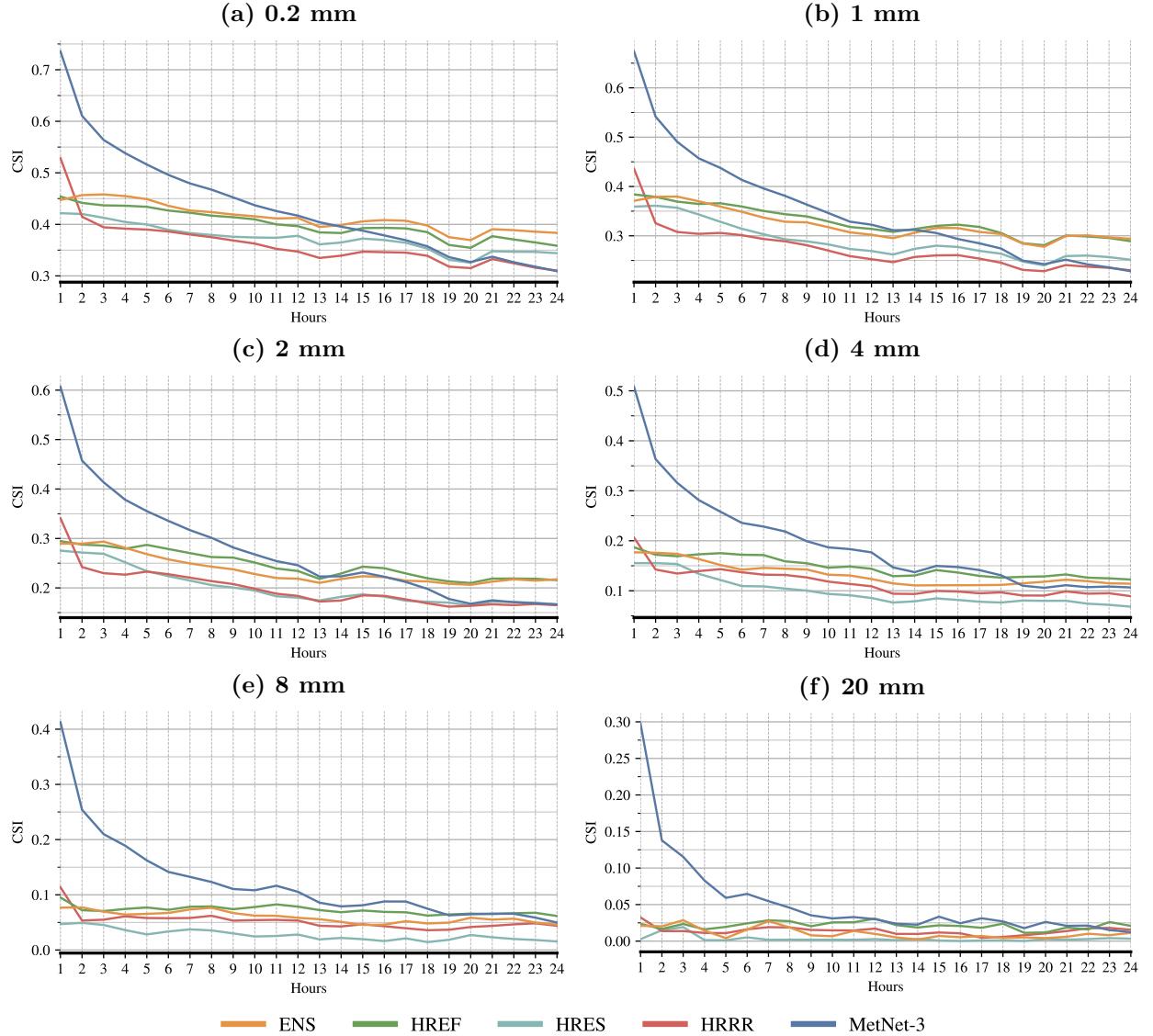
**(a) 0.2 mm**　　**(b) 1 mm**

**(c) 2 mm**　　**(d) 4 mm**

**(e) 8 mm**　　**(f) 20 mm**

ENS　HREF　HRES　HRRR　MetNet-3

Figure 6: CSI values for hourly accumulated precipitation.



**(a) CRPS**　　**(b) MAE**

ENS　HRES　HRRR　MetNet-3　MetNet-3 (hyperlocal)

Figure 7: Performance comparison between MetNet-3 and baselines for U component of wind (i.e. eastward).

**(a) CRPS**

**(b) MAE**

Figure 8: Performance comparison between MetNet-3 and baselines for V component of wind (i.e. north-ward).



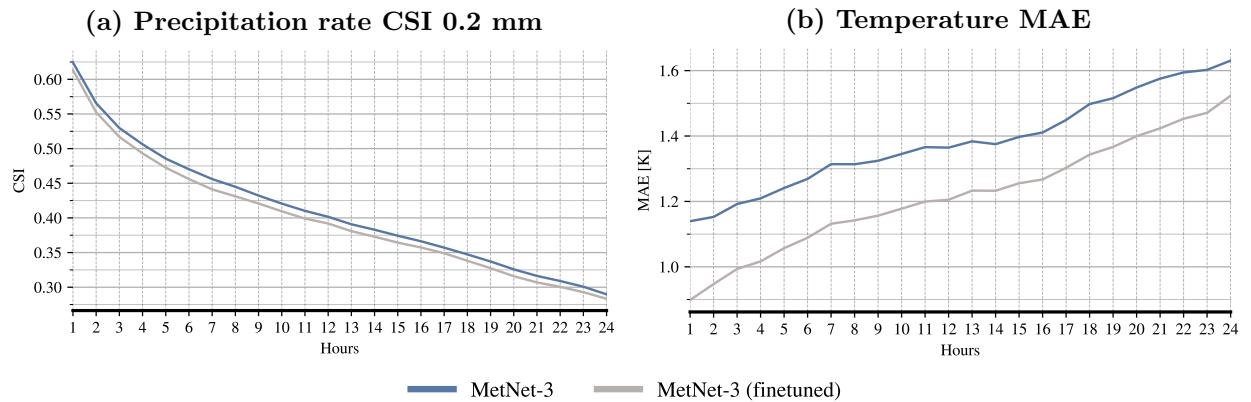**(a) Precipitation rate CSI 0.2 mm**

**(b) Temperature MAE**

Figure 9: Effects of finetuning MetNet-3 for improved OMO performance. Notice that higher CSI and lower MAE are better.
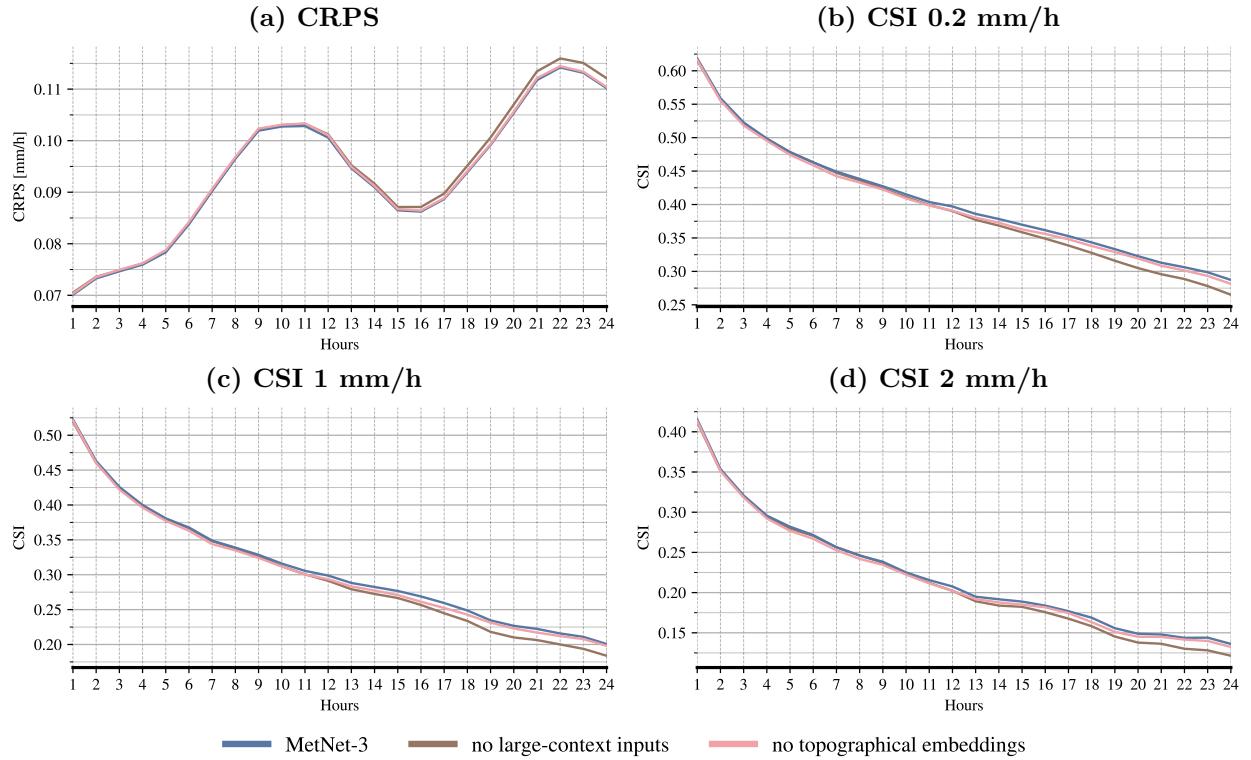
9

Figure 10: Ablations with topographical embeddings and large-context inputs removed on instantaneous precipitation rate. All models shown in this figure have been trained for 150k steps.
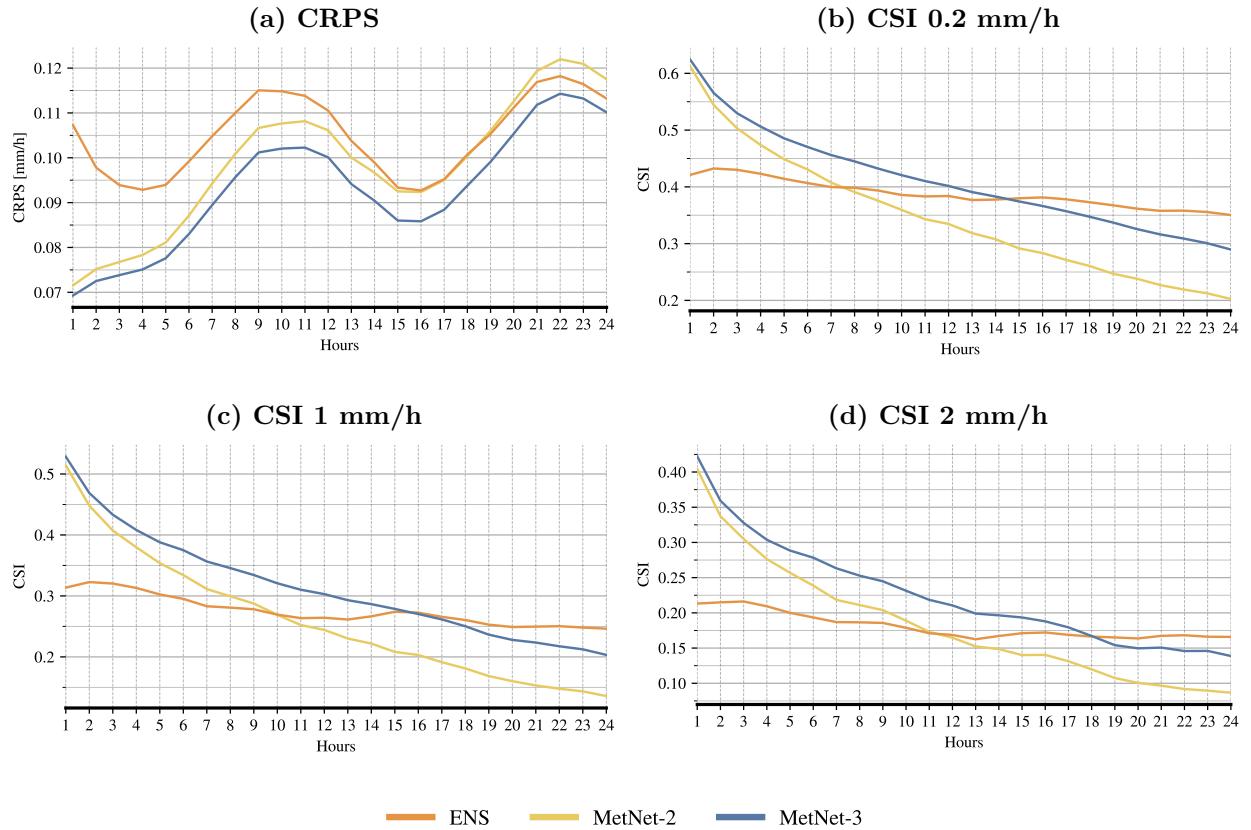
Figure 11: Comparison with MetNet-2 on instantaneous precipitation rate.