

# From Drops to Grid: Noise-Aware Spatio-Temporal Neural Process for Rainfall Estimation

## Supplementary Material

### 6. Derivation of ZIG mean and variance

The model outputs a zero-rain probability  $\pi_0$ , from which we define a deterministic per-sample rain indicator

$$p = \mathbb{1}_{\{1-\pi_0 \geq 0.5\}},$$

so that  $p = 1$  denotes predicted nonzero rainfall and  $p = 0$  otherwise. Given this indicator, the Zero-Inflated Gamma (ZIG) variable  $Y$  is

$$Y \mid (p = 0) = 0, \quad Y \mid (p = 1) \sim \text{Gamma}(\alpha, \beta),$$

with Gamma mean and variance

$$\mu_\Gamma = \frac{\alpha}{\beta}, \quad \sigma_\Gamma^2 = \frac{\alpha}{\beta^2}.$$

Because  $p \in \{0, 1\}$  is deterministic for each sample, conditioning removes all mixture uncertainty:  $Y$  is either identically zero or drawn from a Gamma distribution. Applying the law of total expectation and the law of total variance, respectively:

$$\begin{aligned} \mathbb{E}[Y] &= \mathbb{E}[\mathbb{E}[Y \mid p]] = (1-p) \cdot 0 + p \cdot \mu_\Gamma = p \mu_\Gamma, \\ \text{Var}[Y] &= \mathbb{E}[\text{Var}(Y \mid p)] + \text{Var}(\mathbb{E}[Y \mid p]). \end{aligned}$$

The second term vanishes because  $\mathbb{E}[Y \mid p]$  is constant for the fixed value of  $p$  in a given sample. Therefore,

$$\text{Var}[Y] = p \sigma_\Gamma^2.$$

**Remark.** The more general ZIG variance expression

$$\text{Var}[Y] = p \sigma_\Gamma^2 + p(1-p) \mu_\Gamma^2$$

applies when  $p \in [0, 1]$  represents a *probabilistic mixture weight*, but the deterministic case with indicator  $p$  used here yields the simplified variance  $p \sigma_\Gamma^2$ .

### 7. Training and dataset details

DropsToGrid employs a U-Net of depth 3 with a kernel size of 3, 32 channels, a bottleneck dropout of 0.1, and transformer blocks of depth 2 with 8 heads of dimension 8. The model contains a total of 192K parameters.

It is trained for up to 50K steps with a batch size of 32, and validation is performed every 1,000 steps. Training uses the AdamW optimizer with a cosine-annealed learning rate starting at  $3 \times 10^{-4}$ , weight decay  $1 \times 10^{-4}$ , betas (0.9, 0.999), and an EMA decay of 0.999. On an NVIDIA H100 PCIe GPU, training requires approximately 3 hours.

The dataset spans January-December 2024 and is divided into training, validation, and test periods of 12, 2, and 2 days, respectively, separated by a 12-hour blackout to prevent temporal leakage. This results in 34,287 training, 911 validation, and 914 test patches.

The target region of  $384 \times 384 \text{ km}^2$  is projected onto a 4km/px grid. When multiple stations fall within the same grid cell, the median value is used to mitigate the influence of outliers. Across all samples, up to 902 possible pixels are active in the training set, varying from 797 on January 1st to 884 on December 31st due to station additions, relocations, or temporary outages in the PWS network. The region of interest includes 4,003 land and near-coastal pixels for rainfall estimation. A total of 104 grid cells contain SYNOP stations, with most stations concentrated in urban areas, leading to uneven spatial coverage.

## 8. Metrics and extended results

We evaluate predictions using Critical Success Index (CSI), Fraction Skill Score (FSS), Frequency Bias Index (FBI), and Continuous Ranked Probability Score (CRPS). Additionally, we report Mean Absolute Error (MAE) and Mean Squared Error (MSE), noting that these are sensitive to the prevalence of no-rain events and may be less informative for highly skewed precipitation distributions [3].

**Critical Success Index (CSI)** CSI, or Threat Score, measures event detection accuracy by comparing correctly predicted precipitation events to all predicted or observed events. It balances misses and false alarms, making it especially useful for rare-event forecasting such as heavy rainfall. CSI ranges from 0 (no skill) to 1 (perfect prediction).

$$\text{CSI} = \frac{\text{TP}}{\text{TP} + \text{FP} + \text{FN}}, \quad (1)$$

where TP, FP, and FN are true positives, false positives, and false negatives, respectively, using thresholds of 0.2, 1, 2, 5, and 10 mm/h.

**Fraction Skill Score (FSS)** FSS evaluates the spatial accuracy of predictions by comparing the fraction of predicted and observed positive events within a local neighborhood, rather than pointwise. It accounts for small spatial displacements, making it suitable for high-resolution precipitation estimation. FSS ranges from 0 (no spatial agreement) to 1 (perfect prediction).

$$\text{FSS} = 1 - \frac{\sum_{i=1}^H \sum_{j=1}^W (F_{i,j} - O_{i,j})^2}{\sum_{i=1}^H \sum_{j=1}^W F_{i,j}^2 + \sum_{i=1}^H \sum_{j=1}^W O_{i,j}^2}, \quad (2)$$

where  $F_{i,j}$  and  $O_{i,j}$  are the fractions of predicted and observed positives in the neighborhood of pixel  $(i, j)$ . FSS is computed at thresholds 0.2, 1, 2, 5, and 10 mm/h, and for neighborhood sizes of 2, 10, and 20 pixels.

**Frequency Bias Index (FBI)** FBI quantifies systematic over- or underprediction of events. Values greater than 1 indicate overprediction, while values below 1 indicate underprediction. A bias of 1 indicates that the predicted frequency matches the observed frequency, though not necessarily the spatial or intensity accuracy.

$$\text{FBI} = \frac{\text{TP} + \text{FP}}{\text{TP} + \text{FN}}, \quad (3)$$

evaluated at the same precipitation thresholds as CSI and FSS (0.2, 1, 2, 5, and 10 mm/h).

**Continuous Ranked Probability Score (CRPS)** CRPS evaluates the accuracy of a probabilistic forecast by comparing the predicted cumulative distribution function (CDF) to the observed outcome. Lower values indicate better forecast sharpness, calibration, and reliability, rewarding predictions that assign high probability to the correct rainfall intensity. Computation follows closed-form solutions for Gamma distributions [56]. We report CRPS for the probabilistic models (i.e., ablations and deep learning baselines).

For a Zero-Inflated Gamma (ZIG) forecast, which models rainfall  $R$  as a mixture of a point mass at zero and a continuous Gamma for positive rainfall, CRPS is computed as:

$$\text{CRPS}_{\text{ZIG}} = \begin{cases} p_{\text{nonzero}} \cdot \text{CRPS}_{\Gamma}(y = 0), & R \leq 0.2 \\ p_{\text{zero}} \cdot R + p_{\text{nonzero}} \cdot \text{CRPS}_{\Gamma}(R), & R > 0.2 \end{cases} \quad (4)$$

where  $p_{\text{nonzero}} \in \{0, 1\}$  is the binarized predicted probability of nonzero rainfall,  $p_{\text{zero}} = 1 - p_{\text{nonzero}}$ , and  $\text{CRPS}_{\Gamma}$  is the closed-form CRPS for the Gamma component [56]:

$$\begin{aligned} \text{CRPS}_{\Gamma}(y) = & y(2F_{\Gamma}(y) - 1) - \frac{\alpha}{\beta} (2F_{\Gamma, \alpha+1}(y) - 1) \\ & - \frac{\alpha}{\beta\pi} B\left(\alpha + \frac{1}{2}, \frac{1}{2}\right) \end{aligned} \quad (5)$$

with  $F_{\Gamma}$  the CDF of the Gamma distribution with parameters  $(\alpha, \beta)$ ,  $F_{\Gamma, \alpha+1}$  the CDF with shape  $\alpha + 1$ , and  $B(\cdot, \cdot)$  the Beta function.

**Extended results** Tables 4–9 present a more detailed evaluation of DropsToGrid and the gridded baselines against SYNOP stations, including overall metrics, threshold-specific results for CSI, FBI, and FSS, as well as different neighborhood sizes for FSS.

Table 4. Performance comparison across metrics against operational estimators on research-quality SYNOP stations. DropsToGrid achieves superior performance across all metrics. Best results are shown in **bold**, and second-best in *italics*.

	CSI $\uparrow$	FSS $\uparrow$	FBI $\approx 1$	MAE $\downarrow$	MSE $\downarrow$
OPERA	<i>0.323</i>	<i>0.551</i>	<i>0.688</i>	<i>0.062</i>	<i>0.118</i>
RainViewer	0.304	0.525	2.231	0.089	0.599
IMERG	0.194	0.425	1.315	0.111	0.261
ERA5	0.203	0.363	0.632	0.084	0.174
DropsToGrid	<b>0.551 <math>\pm</math> 0.006</b>	<b>0.795 <math>\pm</math> 0.007</b>	<b>0.995 <math>\pm</math> 0.030</b>	<b>0.037 <math>\pm</math> 0.001</b>	<b>0.070 <math>\pm</math> 0.005</b>

Table 5. Performance comparison across thresholds for CSI against operational estimators on research-quality SYNOP stations. DropsToGrid achieves superior performance across all metrics. Best results are shown in **bold**, and second-best in *italics*.

	CSI $\uparrow$				
	0.2 mm/h	1.0 mm/h	2.0 mm/h	5.0 mm/h	10.0 mm/h
OPERA	<i>0.501</i>	<i>0.460</i>	<i>0.375</i>	<i>0.156</i>	<i>0.125</i>
RadarViewer	0.498	0.446	0.367	0.180	0.029
IMERG	0.280	0.284	0.259	0.102	0.044
ERA5	0.369	0.367	0.281	0.000	0.000
DropsToGrid	<b>0.674 <math>\pm</math> 0.002</b>	<b>0.656 <math>\pm</math> 0.002</b>	<b>0.639 <math>\pm</math> 0.007</b>	<b>0.444 <math>\pm</math> 0.003</b>	<b>0.342 <math>\pm</math> 0.024</b>

Table 6. Performance comparison across thresholds for FBI against operational estimators on research-quality SYNOP stations. DropsToGrid achieves superior performance across all metrics. Best results are shown in **bold**, and second-best in *italics*.

	FBI $\approx 1$				
	0.2 mm/h	1.0 mm/h	2.0 mm/h	5.0 mm/h	10.0 mm/h
OPERA	<i>1.093</i>	0.794	0.657	<i>0.512</i>	0.385
RadarViewer	<b>0.973</b>	<i>1.059</i>	<i>1.076</i>	2.122	5.923
IMERG	1.429	1.459	1.364	1.517	<b>0.808</b>
ERA5	1.346	1.081	0.731	0.000	0.000
DropsToGrid	0.887 $\pm$ 0.005	<b>1.054 <math>\pm</math> 0.005</b>	<b>1.041 <math>\pm</math> 0.020</b>	<b>1.238 <math>\pm</math> 0.038</b>	<i>0.756 <math>\pm</math> 0.101</i>

Table 7. Performance comparison across thresholds for FSS (neighborhood size: 2px) against operational estimators on research-quality SYNOP stations. DropsToGrid achieves superior performance across all metrics. Best results are shown in **bold**, and second-best in *italics*.

	FSS (size: 2px) $\uparrow$				
	0.2 mm/h	1.0 mm/h	2.0 mm/h	5.0 mm/h	10.0 mm/h
OPERA	<i>0.671</i>	<i>0.635</i>	<i>0.550</i>	0.266	<i>0.220</i>
RadarViewer	0.669	0.623	0.543	<i>0.310</i>	0.056
IMERG	0.440	0.444	0.415	0.184	0.084
ERA5	0.542	0.539	0.443	0.000	0.000
DropsToGrid	<b>0.809 <math>\pm</math> 0.001</b>	<b>0.796 <math>\pm</math> 0.002</b>	<b>0.784 <math>\pm</math> 0.005</b>	<b>0.622 <math>\pm</math> 0.003</b>	<b>0.510 <math>\pm</math> 0.031</b>

Table 8. Performance comparison across thresholds for FSS (neighborhood size: 10px) against operational estimators on research-quality SYNOP stations. DropsToGrid achieves superior performance across all metrics. Best results are shown in **bold**, and second-best in *italics*.

	FSS (size: 10px) ↑				
	0.2 mm/h	1.0 mm/h	2.0 mm/h	5.0 mm/h	10.0 mm/h
OPERA	<i>0.780</i>	<i>0.756</i>	0.679	0.338	<i>0.261</i>
RadarViewer	0.777	0.752	<i>0.680</i>	<i>0.419</i>	0.055
IMERG	0.549	0.564	0.535	0.284	0.214
ERA5	0.628	0.640	0.557	0.000	0.000
DropsToGrid	<b>0.900 ± 0.001</b>	<b>0.897 ± 0.001</b>	<b>0.892 ± 0.004</b>	<b>0.782 ± 0.001</b>	<b>0.597 ± 0.039</b>

Table 9. Performance comparison across thresholds for FSS (neighborhood size: 20px) against operational estimators on research-quality SYNOP stations. DropsToGrid achieves superior performance across all metrics. Best results are shown in **bold**, and second-best in *italics*.

	FSS (size: 20px) ↑				
	0.2 mm/h	1.0 mm/h	2.0 mm/h	5.0 mm/h	10.0 mm/h
OPERA	<i>0.843</i>	0.816	0.740	0.418	0.289
RadarViewer	0.839	<i>0.829</i>	<i>0.757</i>	<i>0.488</i>	0.076
IMERG	0.628	0.662	0.636	0.385	<i>0.359</i>
ERA5	0.700	0.730	0.664	0.000	0.000
DropsToGrid	<b>0.943 ± 0.000</b>	<b>0.946 ± 0.001</b>	<b>0.943 ± 0.003</b>	<b>0.860 ± 0.003</b>	<b>0.649 ± 0.040</b>

## 9. Baseline gridded products

For evaluation, we use operational and reanalysis gridded rainfall products as reference baselines rather than ground truth. To ensure comparability, all baselines are resampled to a uniform 4 km grid using bilinear interpolation and converted to hourly rainfall accumulations (mm). Further details on each gridded baseline product are provided below.

**OPERA Odyssey rainfall accumulation.** The EUMETNET OPERA program [29, 53] provides pan-European radar composites combining data from over 160 national radars. The 1-hour accumulation product has 2 km resolution and 15-minute updates. Rain rate is derived from reflectivity using the Marshall-Palmer  $Z = aR^b$  relation ( $a = 200$ ,  $b = 1.6$ ) [40]. Quality filtering employs anomaly removal [47], clutter [57], and beam-blockage [21] corrections. Maximum accumulation is set to 100 mm.

**RainViewer reflectivity.** RainViewer<sup>3</sup> provides near-real time radar composites every 10 minutes at 2km resolution, aggregating over 1,000 stations worldwide. Reflectivity (dBZ) is clipped to [-1, 64] and converted to intensity (mm/h) via Marshall-Palmer [40]. Hourly accumulations are computed by averaging the six 10-min intensity maps.

**IMERG.** The GPM IMERG product [28] merges passive microwave, infrared, and radar data from multiple satellites to provide global rainfall estimates at 0.1° and 30-min intervals. We use the Final Run V07 dataset, bias-corrected to match rain-gauge climatologies. Hourly accumulations are obtained by averaging the two 30-min frames per hour.

**ERA5.** ERA5 [24, 61] is ECMWF’s global reanalysis, combining a frozen numerical model with a consistent data assimilation scheme. Although it does not directly assimilate European rainfall observations, estimates are derived from physically consistent atmospheric fields (humidity, geopotential, winds). ERA5 provides total rainfall, including convective, large-scale, and evaporative effects, at 0.25° resolution and hourly frequency.

**Climate (DMI gridded product).** The DMI *Climate* dataset [55] provides hourly gridded rainfall at 10 km resolution. Fields are generated using Inverse Distance Weighting (IDW) interpolation of quality-controlled gauges, followed by Gaussian smoothing and a coastal-inland correction to adjust distance metrics across land-sea boundaries. The gauge network used to generate the gridded product includes the SYNOP stations previously mentioned.

## 10. Visualizations

Figures 5-8 show visual comparisons of rainfall estimates from DropsToGrid and several operational baselines. DropsToGrid is derived from crowd-sourced PWS stations and RainViewer radar.

The baselines include OPERA radar accumulations, RainViewer reflectivity estimates, IMERG satellite retrievals, ERA5 reanalysis, and DMI’s gridded *Climate* product. CSI is reported for all products against SYNOP stations, except for *Climate*, which is excluded since it incorporates the evaluation SYNOP stations in its gridding.

---

<sup>3</sup><https://www.rainviewer.com>

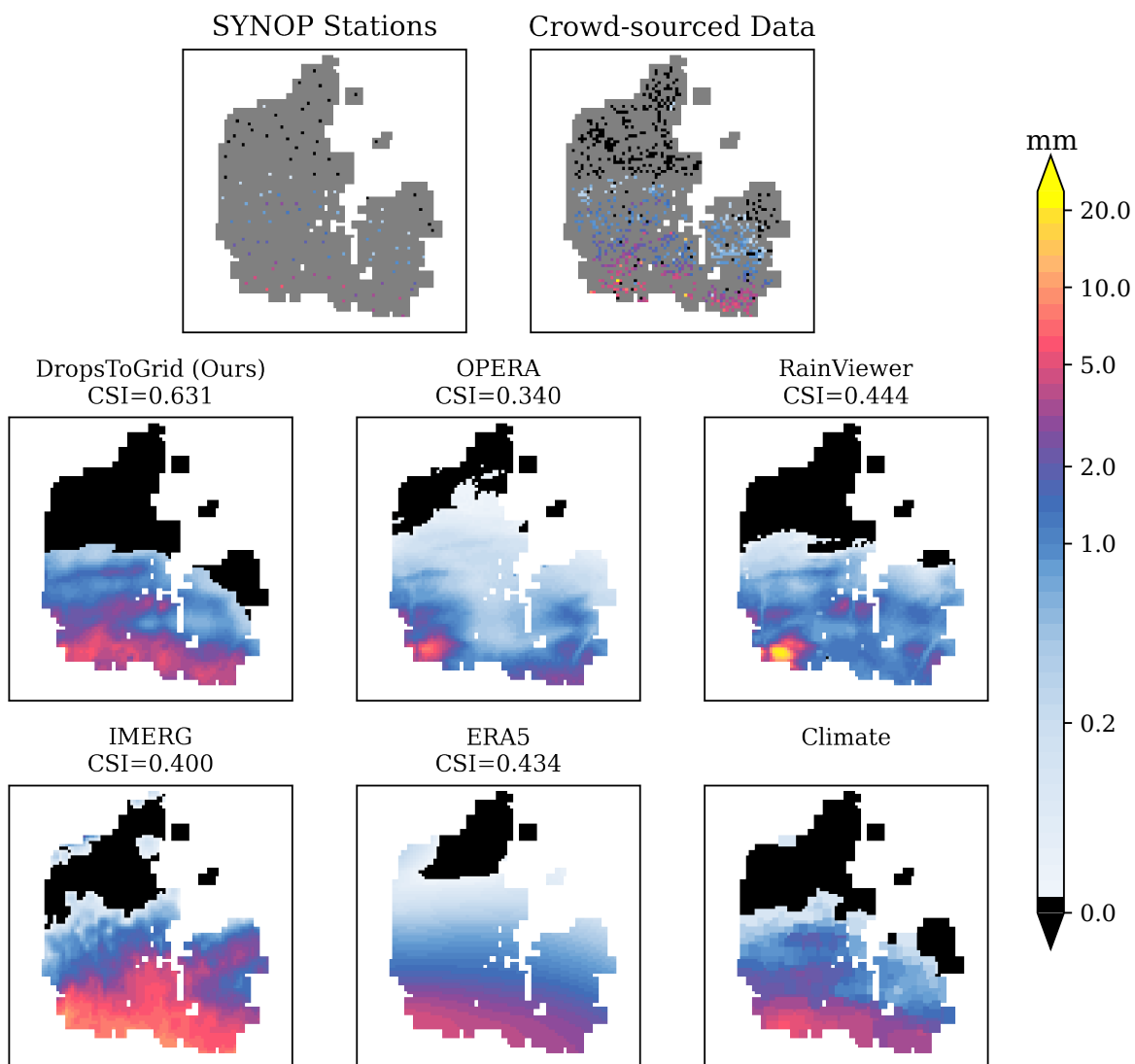


Figure 5. Comparison of rainfall estimators against research-quality SYNOP stations.

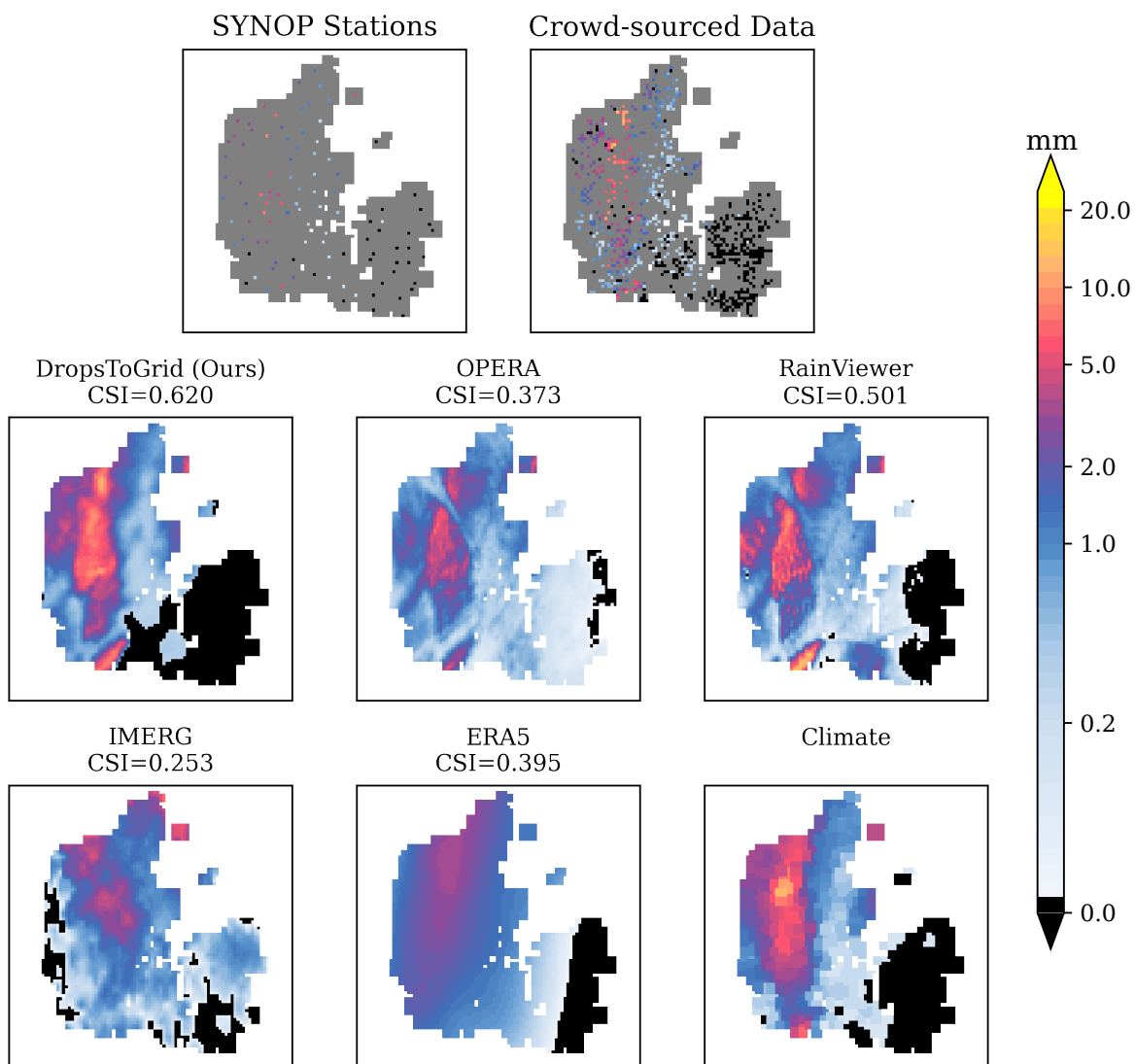


Figure 6. Comparison of rainfall estimators against research-quality SYNOP stations.

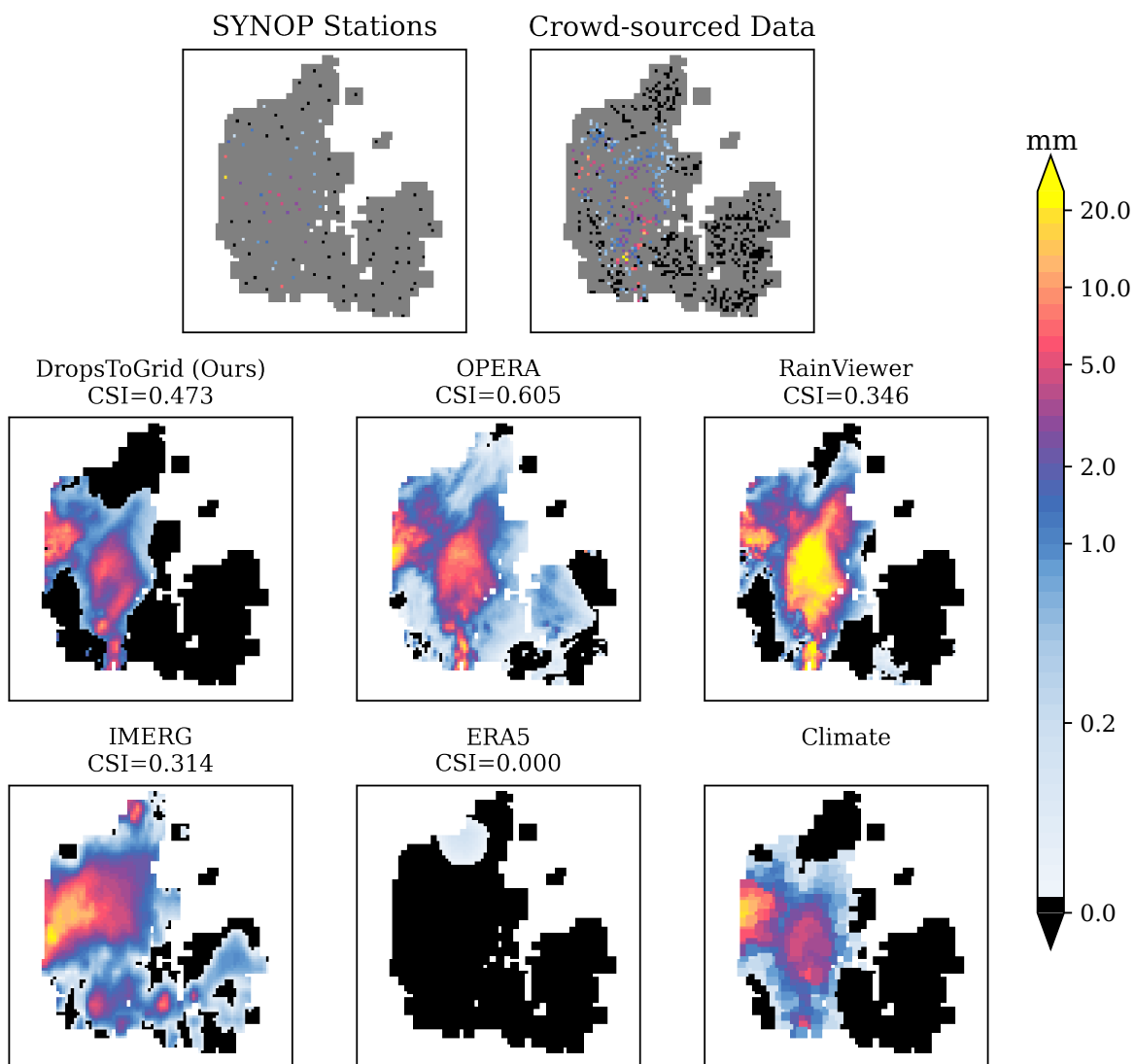


Figure 7. Comparison of rainfall estimators against research-quality SYNOP stations.

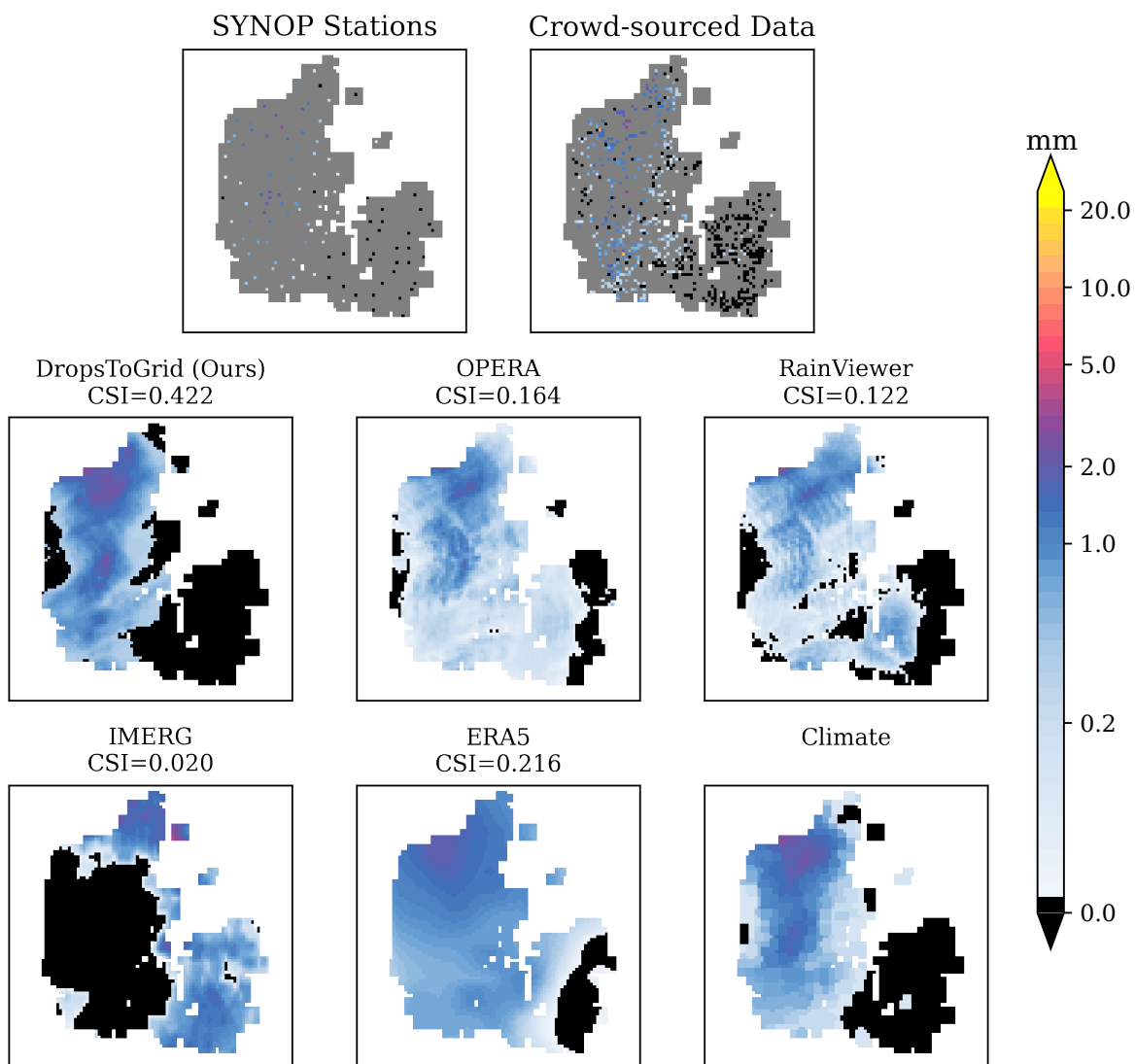


Figure 8. Comparison of rainfall estimators against research-quality SYNOP stations.

## 11. Deep Learning baselines

Beyond the ConvCNP and SwinTNP variants used in the main paper, we evaluate additional baselines. The MM setting uses only station history as input (no radar), while the OOTG setting uses radar and current-time station readings (no history). We further include the translation-equivariant SwinTNP (SwinTNP\_TE) and the approximately translation-equivariant version (SwinTNP\_ATE) from Gridded-TNP [6] in both settings. Finally, we add an extended ALL ConvCNP model that uses both radar and station history to provide a closer comparison to DropsToGrid in terms of data sources.

Across both the PWS holdout and SYNOP stations (Tables 10–11), DropsToGrid achieves the strongest results on almost all metrics, with higher spatial skill (CSI, FSS), better calibration (CRPS), and lower errors (MAE, MSE). ConvCNP models generally outperform the SwinTNP variants in both MM and OOTG settings, likely because they better capture the local structure of rainfall. Even when provided with both radar and station history, the ALL ConvCNP baseline remains below DropsToGrid, underscoring the benefits of DropsToGrid’s processing and fusion design.

Table 10. Performance comparison across metrics against deep learning baselines on PWS holdout stations. DropsToGrid achieves superior performance across most metrics. Best results are shown in **bold**, and second-best in *italics*.

	CSI $\uparrow$	FSS $\uparrow$	CRPS $\downarrow$	FBI $\approx 1$	MAE $\downarrow$	MSE $\downarrow$
<b>MM</b>						
ConvCNP	0.484 $\pm$ 0.012	0.769 $\pm$ 0.035	0.029 $\pm$ 0.000	0.849 $\pm$ 0.061	0.038 $\pm$ 0.001	0.124 $\pm$ 0.003
SwinTNP	0.407 $\pm$ 0.002	0.676 $\pm$ 0.003	0.034 $\pm$ 0.000	0.766 $\pm$ 0.029	0.046 $\pm$ 0.001	0.145 $\pm$ 0.002
SwinTNP_TE	0.412 $\pm$ 0.010	0.672 $\pm$ 0.021	0.034 $\pm$ 0.000	0.814 $\pm$ 0.071	0.047 $\pm$ 0.001	0.147 $\pm$ 0.002
SwinTNP_ATE	0.418 $\pm$ 0.014	0.700 $\pm$ 0.017	0.034 $\pm$ 0.001	0.772 $\pm$ 0.039	0.045 $\pm$ 0.000	0.145 $\pm$ 0.003
<b>OOTG</b>						
ConvCNP	<i>0.486 <math>\pm</math> 0.007</i>	0.748 $\pm$ 0.013	<i>0.027 <math>\pm</math> 0.000</i>	0.752 $\pm$ 0.015	<b>0.035 <math>\pm</math> 0.000</b>	<i>0.115 <math>\pm</math> 0.001</i>
SwinTNP	0.433 $\pm$ 0.001	0.711 $\pm$ 0.006	0.032 $\pm$ 0.000	0.805 $\pm$ 0.020	0.044 $\pm$ 0.000	0.140 $\pm$ 0.002
SwinTNP_TE	0.428 $\pm$ 0.002	0.684 $\pm$ 0.010	0.032 $\pm$ 0.000	0.763 $\pm$ 0.012	0.043 $\pm$ 0.000	0.141 $\pm$ 0.002
SwinTNP_ATE	0.434 $\pm$ 0.004	0.696 $\pm$ 0.010	0.032 $\pm$ 0.000	0.798 $\pm$ 0.045	0.043 $\pm$ 0.001	0.142 $\pm$ 0.003
<b>ALL</b>						
ConvCNP	0.485 $\pm$ 0.022	<i>0.775 <math>\pm</math> 0.035</i>	0.029 $\pm$ 0.000	<b>0.879 <math>\pm</math> 0.109</b>	0.039 $\pm$ 0.001	0.146 $\pm$ 0.022
DropsToGrid	<b>0.532 <math>\pm</math> 0.002</b>	<b>0.819 <math>\pm</math> 0.000</b>	<b>0.026 <math>\pm</math> 0.000</b>	<i>0.877 <math>\pm</math> 0.019</i>	<i>0.035 <math>\pm</math> 0.000</i>	<b>0.112 <math>\pm</math> 0.001</b>

Table 11. Performance comparison across metrics against deep learning baselines on research-quality SYNOP stations. DropsToGrid achieves superior performance across most metrics. Best results are shown in **bold**, and second-best in *italics*.

	CSI $\uparrow$	FSS $\uparrow$	CRPS $\downarrow$	FBI $\approx 1$	MAE $\downarrow$	MSE $\downarrow$
<b>MM</b>						
ConvCNP	0.507 $\pm$ 0.012	<i>0.760 <math>\pm</math> 0.020</i>	0.025 $\pm$ 0.000	0.863 $\pm$ 0.050	0.038 $\pm$ 0.001	0.071 $\pm$ 0.001
SwinTNP	0.424 $\pm$ 0.005	0.666 $\pm$ 0.008	0.030 $\pm$ 0.000	0.797 $\pm$ 0.032	0.045 $\pm$ 0.001	0.093 $\pm$ 0.002
SwinTNP_TE	0.423 $\pm$ 0.014	0.653 $\pm$ 0.027	0.030 $\pm$ 0.000	0.834 $\pm$ 0.089	0.046 $\pm$ 0.001	0.094 $\pm$ 0.001
SwinTNP_ATE	0.437 $\pm$ 0.010	0.679 $\pm$ 0.013	0.030 $\pm$ 0.001	0.780 $\pm$ 0.037	0.045 $\pm$ 0.000	0.091 $\pm$ 0.003
<b>OOTG</b>						
ConvCNP	<i>0.522 <math>\pm</math> 0.015</i>	0.754 $\pm$ 0.021	<i>0.023 <math>\pm</math> 0.000</i>	0.823 $\pm$ 0.031	<b>0.035 <math>\pm</math> 0.000</b>	<b>0.067 <math>\pm</math> 0.000</b>
SwinTNP	0.457 $\pm$ 0.005	0.700 $\pm$ 0.007	0.028 $\pm$ 0.000	0.859 $\pm$ 0.035	0.043 $\pm$ 0.001	0.087 $\pm$ 0.001
SwinTNP_TE	0.439 $\pm$ 0.001	0.655 $\pm$ 0.005	0.028 $\pm$ 0.000	0.790 $\pm$ 0.009	0.042 $\pm$ 0.000	0.087 $\pm$ 0.001
SwinTNP_ATE	0.439 $\pm$ 0.007	0.665 $\pm$ 0.014	0.028 $\pm$ 0.000	0.824 $\pm$ 0.045	0.043 $\pm$ 0.001	0.087 $\pm$ 0.003
<b>ALL</b>						
ConvCNP	0.513 $\pm$ 0.007	0.757 $\pm$ 0.006	0.025 $\pm$ 0.000	<i>0.987 <math>\pm</math> 0.144</i>	0.039 $\pm$ 0.001	0.090 $\pm$ 0.016
DropsToGrid	<b>0.551 <math>\pm</math> 0.006</b>	<b>0.795 <math>\pm</math> 0.007</b>	<b>0.023 <math>\pm</math> 0.000</b>	<b>0.995 <math>\pm</math> 0.030</b>	<i>0.037 <math>\pm</math> 0.001</i>	<i>0.070 <math>\pm</math> 0.005</i>

## 12. Ablation

Tables 12 and 13 summarize the ablation studies conducted to assess the contribution of each component of DropsToGrid on the PWS holdout stations and the research-grade SYNOP stations, respectively. The two primary ablations discussed in the main paper examine (i) replacing the carefully designed fusion bottleneck with a standard convolution that simply stacks all latent source representations in the input (*no\_bottleneck*), and (ii) allowing input stations to also serve as target stations during training (*target\_inputs*), in contrast to DropsToGrid’s strategy of excluding inputs from the prediction targets to avoid direct input-output mapping and mitigate noisy observations. We additionally evaluate variants without station history (*no\_stations*), without radar (*no\_radar*), and with a transformer using standard (non-translation-equivariant) attention (*no\_te*). Regarding the output distribution, beyond the *target\_inputs* ablation, we test a plain gamma distribution that omits the zero-inflation component (*gamma*), as well as a Gaussian output distribution (*gaussian*), to highlight the importance of modeling the highly skewed, zero-dominated nature of rainfall.

As noted in the main paper, the specialized fusion bottleneck is critical, with *no\_bottleneck* showing consistent degradation across all metrics. Although the effects are smaller in magnitude, radar information, station history, and translation-equivariant attention each provide measurable gains, from pixelwise skill (CSI) to calibration of the distribution (CRPS) and overall error (MAE/MSE). Preventing direct input–output mapping is also important, particularly given the noise in PWS observations, as reflected in the significant performance drop against SYNOP stations in *target\_inputs*. Finally, the *gamma* and *gaussian* variants confirm that an appropriate output distribution is essential for rainfall, whose statistics are highly skewed and dominated by zero-rain cases.

Table 12. Ablation experiments on PWS holdout stations. DropsToGrid achieves superior performance across most metrics. Best results are shown in **bold**, and second-best in *italics*.

	CSI $\uparrow$	FSS $\uparrow$	CRPS $\downarrow$	FBI $\approx 1$	MAE $\downarrow$	MSE $\downarrow$
DropsToGrid	<b>0.532 <math>\pm</math> 0.002</b>	0.819 $\pm$ 0.000	<b>0.026 <math>\pm</math> 0.000</b>	0.877 $\pm$ 0.019	<b>0.035 <math>\pm</math> 0.000</b>	<b>0.112 <math>\pm</math> 0.001</b>
no_bottleneck	0.520 $\pm$ 0.015	0.804 $\pm$ 0.025	0.027 $\pm$ 0.000	0.868 $\pm$ 0.042	0.036 $\pm$ 0.001	0.115 $\pm$ 0.005
no_stations	<i>0.530 <math>\pm</math> 0.004</i>	0.825 $\pm$ 0.006	<i>0.026 <math>\pm</math> 0.000</i>	0.905 $\pm$ 0.006	0.036 $\pm$ 0.000	0.113 $\pm$ 0.003
no_radar	0.530 $\pm$ 0.004	<b>0.829 <math>\pm</math> 0.004</b>	0.027 $\pm$ 0.000	<i>0.905 <math>\pm</math> 0.011</i>	0.036 $\pm$ 0.000	0.115 $\pm$ 0.001
no_te	0.530 $\pm$ 0.005	<i>0.826 <math>\pm</math> 0.005</i>	0.026 $\pm$ 0.000	0.871 $\pm$ 0.018	<i>0.036 <math>\pm</math> 0.000</i>	<i>0.113 <math>\pm</math> 0.001</i>
target_inputs	0.514 $\pm$ 0.008	0.816 $\pm$ 0.006	0.028 $\pm$ 0.001	<b>0.996 <math>\pm</math> 0.044</b>	0.039 $\pm$ 0.001	0.149 $\pm$ 0.007
gamma	0.471 $\pm$ 0.016	0.727 $\pm$ 0.024	0.032 $\pm$ 0.000	0.705 $\pm$ 0.052	0.038 $\pm$ 0.001	0.118 $\pm$ 0.003
gaussian	0.414 $\pm$ 0.006	0.637 $\pm$ 0.016	0.045 $\pm$ 0.003	0.635 $\pm$ 0.027	0.041 $\pm$ 0.001	0.129 $\pm$ 0.001

Table 13. Ablation experiments on research-quality SYNOP stations. DropsToGrid achieves superior performance across most metrics. Best results are shown in **bold**, and second-best in *italics*.

	CSI $\uparrow$	FSS $\uparrow$	CRPS $\downarrow$	FBI $\approx 1$	MAE $\downarrow$	MSE $\downarrow$
DropsToGrid	<b>0.551 <math>\pm</math> 0.006</b>	<i>0.795 <math>\pm</math> 0.007</i>	<b>0.023 <math>\pm</math> 0.000</b>	<b>0.995 <math>\pm</math> 0.030</b>	<i>0.037 <math>\pm</math> 0.001</i>	<b>0.070 <math>\pm</math> 0.005</b>
no_bottleneck	0.529 $\pm$ 0.023	0.770 $\pm$ 0.023	<i>0.024 <math>\pm</math> 0.000</i>	0.966 $\pm$ 0.058	<b>0.037 <math>\pm</math> 0.000</b>	0.077 $\pm$ 0.006
no_stations	0.549 $\pm$ 0.003	0.794 $\pm$ 0.004	0.024 $\pm$ 0.000	1.032 $\pm$ 0.019	0.037 $\pm$ 0.000	0.078 $\pm$ 0.001
no_radar	<i>0.551 <math>\pm</math> 0.009</i>	<b>0.796 <math>\pm</math> 0.009</b>	0.024 $\pm$ 0.000	1.026 $\pm$ 0.025	0.037 $\pm$ 0.000	0.072 $\pm$ 0.002
no_te	0.542 $\pm$ 0.002	0.789 $\pm$ 0.004	0.024 $\pm$ 0.000	<i>0.989 <math>\pm</math> 0.006</i>	0.037 $\pm$ 0.000	<i>0.072 <math>\pm</math> 0.003</i>
target_inputs	0.505 $\pm$ 0.009	0.755 $\pm$ 0.008	0.029 $\pm$ 0.001	1.198 $\pm$ 0.044	0.042 $\pm$ 0.001	0.119 $\pm$ 0.014
gamma	0.485 $\pm$ 0.014	0.719 $\pm$ 0.022	0.033 $\pm$ 0.000	0.769 $\pm$ 0.041	0.038 $\pm$ 0.000	0.075 $\pm$ 0.003
gaussian	0.409 $\pm$ 0.009	0.600 $\pm$ 0.016	0.046 $\pm$ 0.002	0.622 $\pm$ 0.028	0.041 $\pm$ 0.001	0.081 $\pm$ 0.002

### 13. Station analysis

To assess how varying observational coverage affects DropsToGrid, we study performance under progressively reduced densities of input PWS stations. Starting from all 902 pixels with PWS data, we randomly mask stations in 10% increments, using nested masks so that each higher-density configuration contains all stations from the previous one. We repeat the experiment with three seeds and average the results, where each seed influences the station masking and model. For each density level, we evaluate the complete test set against SYNOP stations using CRPS and CSI with only the corresponding subset of stations as input.

DropsToGrid performs best when all stations are available, but remains highly robust even under severe sparsity: at 5% density, performance remains strong and still surpasses operational estimators such as OPERA (CSI 0.323) by more than 24% in CSI. This demonstrates the model’s ability to leverage limited point observations while still capturing localized rainfall structure. Figure 9 summarizes the performance trend, and Figure 10 illustrates a qualitative example showing the station inputs (5%, 30%, and 100%), radar observations, the resulting predictions by DropsToGrid, and SYNOP targets, along with the corresponding CSI values. As shown in the case study, with more input stations, the level of detail of DropsToGrid increases.

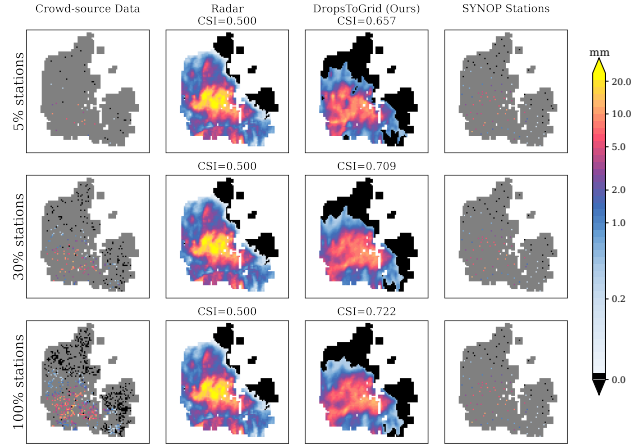


Figure 10. Qualitative example of station masked inputs, radar input, predictions, and targets.

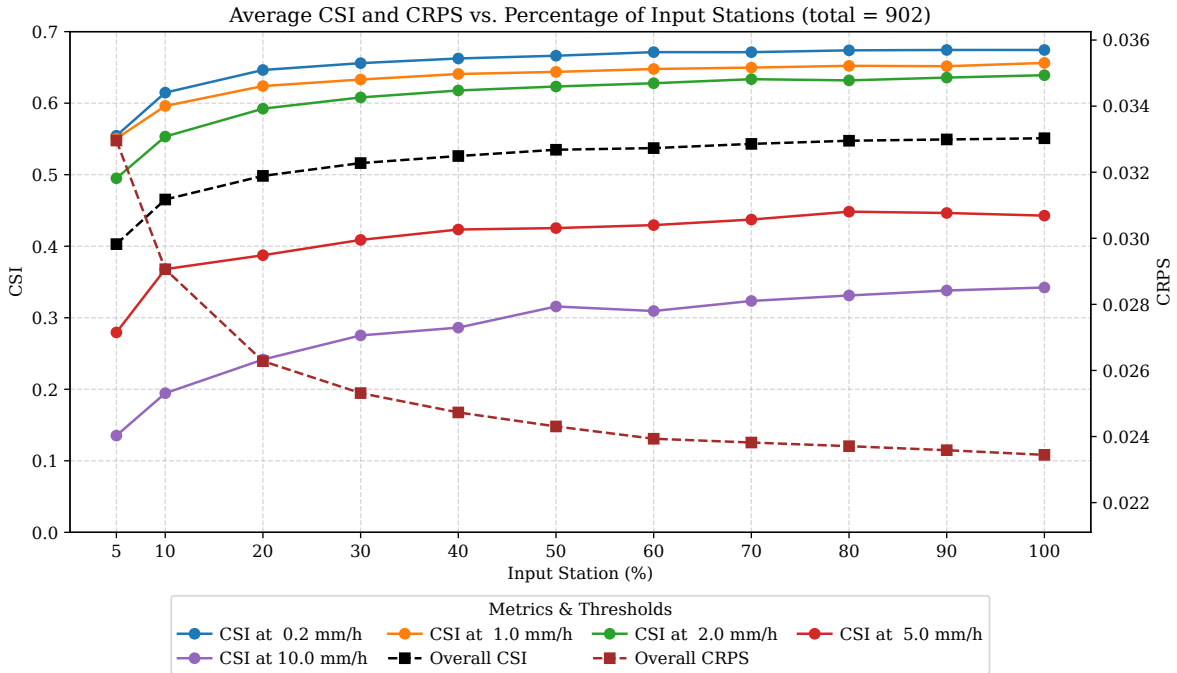


Figure 9. Average performance across station density levels.

## 14. Europe-wide densification

We evaluate DropsToGrid across the entirety of Europe for the full year of 2025, a period entirely unseen during training, encompassing a wide range of climatic and topographic conditions. This evaluation requires multiple sources: the WeatherUnderground PWS Network<sup>4</sup>, OPERA radar, and GHCNh SYNOP stations<sup>5</sup>, all providing at least EU-wide coverage. Table 14 presents results in which a model trained exclusively on Danish data (DropsToGrid\_dk) is evaluated across Europe against both operational baselines and a learned baseline (ConvCNP\_dk).

We further demonstrate scalability and improved performance when the same compact, parameter-efficient model is trained on all available EU data (DropsToGrid\_eu). While training still uses patch-based inputs, inference in DropsToGrid is patch-free, achieved through matched TE-Transformer windows, which eliminates border artifacts. EU-wide inference is efficient, taking 0.20s for DropsToGrid versus 0.47s for ConvCNP, the latter being slower due to patch-based processing and 3D convolutions.

Climate baseline is excluded from EU-wide evaluation because it is only available for Denmark and utilizes SYNOP stations. Consequently, Table 14 reflects the limitations of global PWS and SYNOP networks, including lower data quality, sparser coverage, and heterogeneous temporal sampling compared to the Denmark-focused dataset (Table 1), leading to overall performance drops across baselines and learned models. Despite these challenges, DropsToGrid consistently outperforms all baselines. In contrast, ConvCNP trained in Denmark suffers from missing temporal attention and cross-attention mechanisms, limiting its

ability to handle longer temporal gaps and heterogeneous spatiotemporal sources.

Overall, DropsToGrid exhibits robust generalization to unseen regions and seasons across Europe in 2025. It learns to suppress noise from PWS temporal observations, rely on radar in sparsely observed areas, and correct biases in radar-only regions by capturing systematic noise patterns. These capabilities allow DropsToGrid to substantially surpass operational baselines, regardless of network quality, coverage, or temporal heterogeneity.

<sup>4</sup><https://www.wunderground.com/pws/overview>

<sup>5</sup><https://www.ncei.noaa.gov/products/global-historical-climatology-network-hourly>

Table 14. Performance comparison across regions and station types for all of 2025. **Note:** PWS and SYNOP stations differ from those in the main paper, going from dense DK to sparser and noisier but wider EU coverage. With worse data quality and access, overall performance is lower, yet, DropsToGrid maintains similar improvements over baselines as in Table 1, demonstrating robust rainfall estimation.

Method	Denmark						Europe					
	PWS holdout			SYNOP Stations			PWS holdout			SYNOP Stations		
	CSI↑	FSS↑	MAE↓	CSI↑	FSS↑	MAE↓	CSI↑	FSS↑	MAE↓	CSI↑	FSS↑	MAE↓
OPERA	0.206	0.395	0.068	0.281	0.409	0.123	0.298	0.608	0.082	0.261	0.455	0.093
RainViewer	0.240	0.471	0.069	0.289	0.430	0.130	0.314	0.619	0.086	0.308	0.540	0.113
IMERG	0.182	0.392	0.121	0.232	0.351	0.216	0.199	0.473	0.120	0.205	0.426	0.148
ERA5	0.157	0.350	0.103	0.197	0.299	0.161	0.144	0.321	0.119	0.171	0.324	0.137
Climate	0.388	0.713	0.057	-	-	-	-	-	-	-	-	-
ConvCNP_dk	0.330	0.624	0.059	0.337	0.462	0.123	0.253	0.466	0.085	0.243	0.407	0.110
DropsToGrid_dk	<b>0.451</b>	<b>0.741</b>	<b>0.046</b>	0.465	0.606	<b>0.093</b>	0.462	0.785	0.050	0.387	0.628	<b>0.079</b>
DropsToGrid_eu	0.444	0.729	<b>0.046</b>	<b>0.472</b>	<b>0.614</b>	0.096	<b>0.476</b>	<b>0.802</b>	<b>0.049</b>	<b>0.400</b>	<b>0.646</b>	<b>0.079</b>

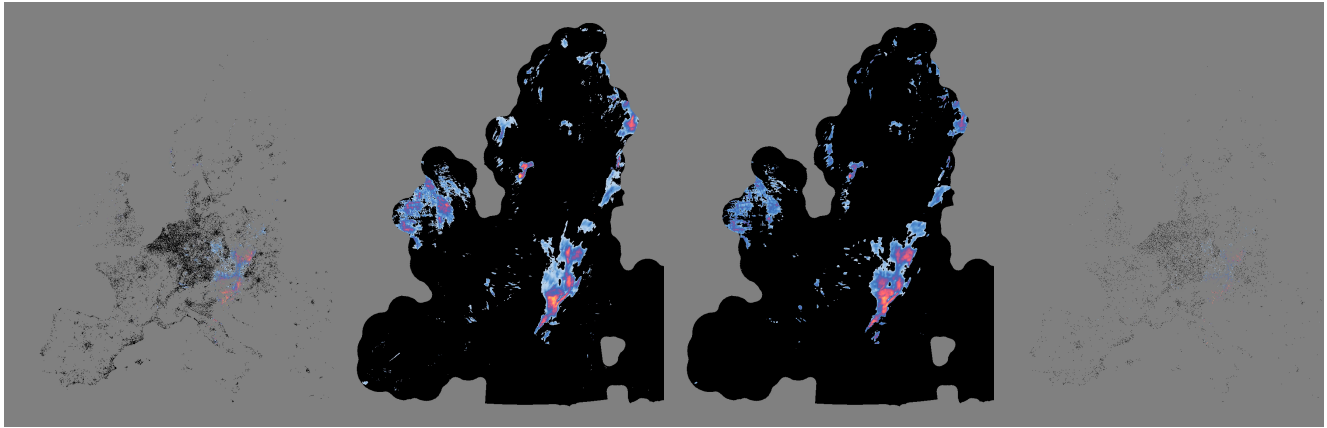


Figure 11. Sample EU densification. From left to right: input PWS stations, input radar, DropsToGrid prediction, and PWS holdout stations.

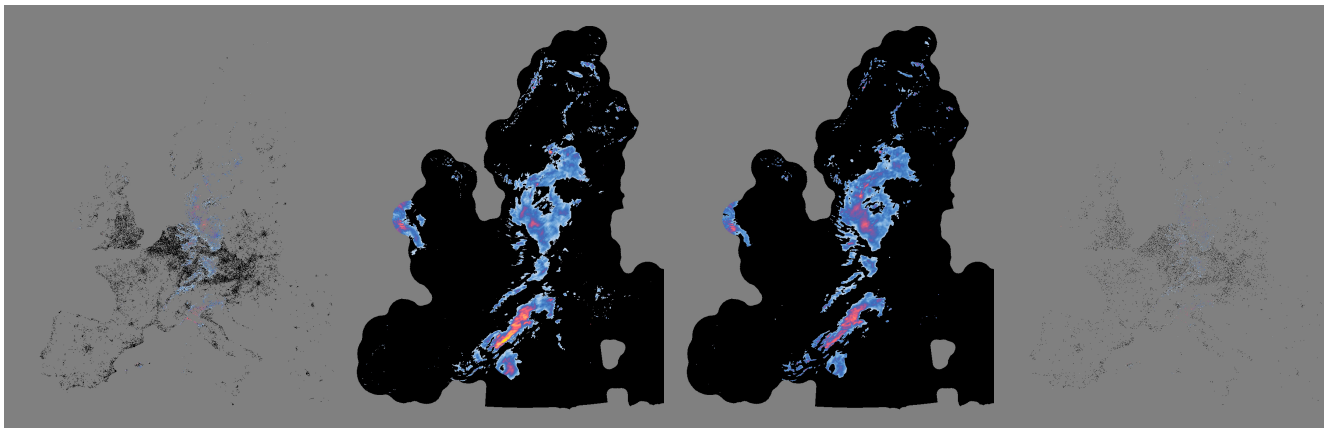


Figure 12. Sample EU densification. From left to right: input PWS stations, input radar, DropsToGrid prediction, and PWS holdout stations.