

Supporting Information for “Reliable precipitation nowcasting using probabilistic diffusion model”

Congyi Nai^{1,3}, Baoxiang Pan², Jiarui Hai⁴, Xi Chen², Qihong Tang^{1,3}, Guangheng Ni⁴,
Qingyun Duan⁵, Bo Lu⁶, Ziniu Xiao², Xingcai Liu^{1,3,*}

¹ Key Laboratory of Water Cycle and Related Land Surface Processes, Institute of Geographic
Sciences and Natural Resources Research, Chinese Academy of Sciences, Beijing, China

² Institute of Atmospheric physics, Chinese Academy of Sciences, Beijing, China

³ University of Chinese Academy of Sciences, Beijing, China

⁴ State Key Laboratory of Hydro-science and Engineering, Department of Hydraulic
Engineering, Tsinghua University, Beijing 100084, China

⁵ The National Key Laboratory of Water Disaster Prevention, Hohai University, Nanjing,
China

⁶ Laboratory for Climate Studies and CMA-NJU Joint Laboratory for Climate Prediction
Studies, National Climate Center, China Meteorological Administration, Beijing, China

*Corresponding author. E-mail address: xingcailiu@igsnnr.ac.cn

1 Details of diffusion model

1.1 Basic diffusion

Let x_0 be a sample from the data distribution $q(x_0)$, and defines a sequence of
increasingly noisy versions of x which we call the latent variables x_t ($t = 1 \dots T$)
through the forward diffusion process, described by

$$q(x_t|x_{t-1}) = N(x_t; \sqrt{1 - \beta_t}x_{t-1}, \beta_t I) \quad (1)$$

Then, the form of $q(x_t|x_0)$ can be recursively derived through repeated
applications of the reparameterization trick, suppose we have $\{\epsilon_t, \bar{\epsilon}_t\}_{t=0}^T \sim N(0, I)$, Then, for
an arbitrary sample $x_t \sim q(x_t|x_0)$, we can rewrite it as:

$$\begin{aligned} x_t &= \sqrt{1 - \beta_t}x_{t-1} + \sqrt{\beta_t}\epsilon_{t-1} \\ &= \sqrt{\alpha_t}(\sqrt{\alpha_{t-1}}x_{t-2} + \sqrt{1 - \alpha_{t-1}}\epsilon_{t-2}) + \sqrt{1 - \alpha_t}\epsilon_{t-1} \\ &= \sqrt{\alpha_t\alpha_{t-1}}x_{t-2} + \sqrt{\alpha_t - \alpha_t\alpha_{t-1}}\epsilon_{t-2} + \sqrt{1 - \alpha_t}\epsilon_{t-1} \\ &= \sqrt{\alpha_t\alpha_{t-1}}x_{t-2} + \sqrt{1 - \alpha_t\alpha_{t-1}}\bar{\epsilon}_{t-2} \\ &= \dots \\ &= \sqrt{\prod_{i=1}^t \alpha_i}x_0 + \sqrt{1 - \prod_{i=1}^t \alpha_i}\bar{\epsilon}_0 \\ &= \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\bar{\epsilon}_0, \text{ where } \bar{\alpha}_t = \prod_{i=1}^t \alpha_i \end{aligned} \quad (2)$$

In equation 3, we have leveraged the property that the sum of two independent
Gaussian random variables retains a Gaussian distribution, with the mean being the sum
of the two individual means and the variance being the sum of their variances.
 $\sqrt{\alpha_t - \alpha_t\alpha_{t-1}}\epsilon_{t-2}$ is a sample from Gaussian $N(0, (\alpha_t - \alpha_t\alpha_{t-1})I)$, $\sqrt{1 - \alpha_t}\epsilon_{t-1}$ is a sample
from Gaussian $N(0, (1 - \alpha_t)I)$, we can then treat their sum as a random variable
sampled from Gaussian $N(0, (1 - \alpha_t + \alpha_t - \alpha_t\alpha_{t-1})I)$. Hence, the x_t can be sampled
directly from x_0 , the transition kernel is

$$q(x_t|x_0) = N(x_t; \sqrt{\bar{\alpha}_t}x_0, \sqrt{1 - \bar{\alpha}_t}I), \text{ where } \alpha_t = 1 - \beta_t, \bar{\alpha}_t = \prod_{i=1}^t \alpha_i. \quad (3)$$

38 Given X_0 and a Gaussian vector $\epsilon \sim N(0, I)$ and applying the transformation
 $X_t = \sqrt{\bar{\alpha}_t} X_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon.$ (4)

39 When the $\bar{\alpha}_T \rightarrow 0$, X_T is well approximated by Gaussian distribution. During the
 40 forward process, noise is gradually added to the data until it loses its original spatial
 41 structure characteristics and becomes pure noise. If we can solve the reverse process
 42 $P(X_{t-1}|X_t)$, we can sample $X_T \sim N(0, I)$ then a sequence of neural networks is employed to
 43 gradually reduce the noise in a series of steps $X_T, X_{T-1} \dots X_0$. These properties suggest
 44 learning a learnable Markov chain model $P_\theta(X_{t-1}|X_t)$ to approximate the true reverse
 45 process:

$$P_\theta(X_{t-1}|X_t) = N(X_{t-1}; \mu_\theta(X_t), \Sigma_\theta(X_t)), \quad (5)$$

46 Therefore, in a diffusion model, we are only interested in learning conditionals
 47 $P_\theta(X_{t-1}|X_t)$, the diffusion model can be optimized by maximizing the variational lower
 48 bound (VLB) of the log-likelihood of the data X_0 ,

$$\begin{aligned} 49 \quad \mathbb{E}_{q(X_0)}(-\log P_\theta(X_0)) &\leq \mathbb{E}_{q(X_0)}[-\log P_\theta(X_0) + D_{KL}(q(X_{1:T}|X_0) || P_\theta(X_{1:T}|X_0))] \\ 50 &= \mathbb{E}_{q(X_0)} \left[-\log P_\theta(X_0) + \int q(X_{1:T}|X_0) \log \frac{q(X_{1:T}|X_0)}{P_\theta(X_{0:T})/P_\theta(X_0)} dX_{1:T} \right] \\ 51 &= \mathbb{E}_{q(X_0)} \left[-\log P_\theta(X_0) + \int q(X_{1:T}|X_0) \log \frac{q(X_{1:T}|X_0)}{P_\theta(X_{0:T})} dX_{1:T} + \log P_\theta(X_0) \right] \\ &= \mathbb{E}_{q(X_{0:T})} \log \frac{q(X_{1:T}|X_0)}{P_\theta(X_{0:T})} = L_{VLB} \end{aligned} \quad (6)$$

52 We can rewrite variational lower bound (VLB) as,

$$\begin{aligned} 53 \quad L_{VLB} &= \mathbb{E}_{q(x_0:T)} [\log \frac{q(x_{1:T}|x_0)}{p_\theta(x_{0:T})}] \\ 54 &= \mathbb{E}_q [\log \frac{\prod_{t=1}^T q(x_t|x_{t-1})}{p_\theta(x_T) \prod_{t=1}^T p_\theta(x_{t-1}|x_t)}] \\ 55 &= \mathbb{E}_q [-\log p_\theta(x_T) + \sum_{t=1}^T \log \frac{q(x_t|x_{t-1})}{p_\theta(x_{t-1}|x_t)}] \\ 56 &= \mathbb{E}_q [-\log p_\theta(x_T) + \sum_{t=2}^T \log \frac{q(x_t|x_{t-1})}{p_\theta(x_{t-1}|x_t)} + \log \frac{q(x_1|x_0)}{p_\theta(x_0|x_1)}] \\ 57 &= \mathbb{E}_q [-\log p_\theta(x_T) + \sum_{t=2}^T \log \left(\frac{q(x_{t-1}|x_t, x_0)}{p_\theta(x_{t-1}|x_t)} \cdot \frac{q(x_t|x_0)}{q(x_{t-1}|x_0)} \right) + \log \frac{q(x_1|x_0)}{p_\theta(x_0|x_1)}] \\ 58 &= \mathbb{E}_q [-\log p_\theta(x_T) + \sum_{t=2}^T \log \frac{q(x_{t-1}|x_t, x_0)}{p_\theta(x_{t-1}|x_t)} + \sum_{t=2}^T \log \frac{q(x_t|x_0)}{q(x_{t-1}|x_0)} + \log \frac{q(x_1|x_0)}{p_\theta(x_0|x_1)}] \\ 59 &= \mathbb{E}_q [-\log p_\theta(x_T) + \sum_{t=2}^T \log \frac{q(x_{t-1}|x_t, x_0)}{p_\theta(x_{t-1}|x_t)} + \log \frac{q(x_T|x_0)}{q(x_1|x_0)} + \log \frac{q(x_1|x_0)}{p_\theta(x_0|x_1)}] \\ 60 &= \mathbb{E}_q [\log \frac{q(x_T|x_0)}{p_\theta(x_T)} + \sum_{t=2}^T \log \frac{q(x_{t-1}|x_t, x_0)}{p_\theta(x_{t-1}|x_t)} - \log p_\theta(X_0|X_1)] \\ &= \mathbb{E}_q [D_{KL}(q(X_T|X_0) || p_\theta(X_T)) + \sum_{t=2}^T D_{KL}(q(X_{t-1}|X_t, X_0) || p_\theta(X_{t-1}|X_t) - \log p_\theta(X_0|X_1))] \end{aligned} \quad (7)$$

61 This formulation also has an elegant interpretation, which is revealed when
 62 inspecting each individual term:

- 63 1. $L_0 = \mathbb{E}_q[\log p_\theta(X_0|X_1)]$ can be interpreted as a reconstruction term.
- 64 2. $L_T = \mathbb{E}_q[D_{KL}(q(X_T|X_0) || p_\theta(X_T))]$ represents how close the distribution of the final
 65 noisified input is to the standard Gaussian prior, is equal to zero under our

assumptions.

3. $L_t = \mathbb{E}_q[\sum_{t=2}^T D_{KL}(q(X_{t-1}|X_t, X_0)||p_\theta(X_{t-1}|X_t))]$ is a denoising matching term. The $q(X_{t-1}|X_t, X_0)$ acts as a ground-truth signal and $p_\theta(X_{t-1}|X_t)$ is our desired denoising transition step. This term is therefore minimized when the two denoising steps match as closely as possible. It is the primary optimization objective.

If we have knowledge of X_0 , we can obtain $q(X_{t-1}|X_t, X_0)$ through the Bayes' theorem,

$$\begin{aligned} q(X_{t-1}|X_t, X_0) &= q(X_t|X_{t-1}, X_0) \frac{q(X_{t-1}|X_0)}{q(X_t|X_0)} \\ &\propto \exp\left(-\frac{1}{2}\left(\frac{(X_t - \sqrt{\alpha_t}X_{t-1})^2}{\beta_t} + \frac{(X_{t-1} - \sqrt{\alpha_{t-1}}X_0)^2}{1 - \alpha_{t-1}} - \frac{(X_t - \sqrt{\alpha_t}X_0)^2}{1 - \alpha_t}\right)\right) \\ &= \exp\left(-\frac{1}{2}\left(\left(\frac{\alpha_t}{\beta_t} + \frac{1}{1 - \alpha_{t-1}}\right)X_{t-1}^2 + \left(\frac{2\sqrt{\alpha_t}}{\beta_t}X_t + \frac{2\sqrt{\alpha_{t-1}}}{1 - \alpha_{t-1}}X_0\right)X_{t-1} + C(X_t, X_0)\right)\right) \\ &= N(X_{t-1}; \tilde{\mu}(X_t, X_0), \tilde{\beta}(t)I) \end{aligned} \quad (8)$$

Recall equation 8 and equation 5, we can obtain,

$$\tilde{\mu}_\theta(X_t, X_0) = \frac{1}{\sqrt{\alpha_t}}\left(X_t - \frac{1 - \alpha_t}{\sqrt{1 - \alpha_t}}\epsilon_\theta(X_t, t)\right) \quad (9)$$

Let us consider the $L_t = D_{KL}(q(X_{t-1}|X_t, X_0)||p_\theta(X_{t-1}|X_t))$, given equation 6 and equation 8, we can get the loss function,

$$\begin{aligned} L_t &= \mathbb{E}_{x_0, \epsilon} \left[\frac{1}{2\|\Sigma_\theta(x_t, t)\|_2^2} \|\tilde{\mu}_t(x_t, x_0) - \mu_\theta(x_t, t)\|_2^2 \right] \\ &= \mathbb{E}_{x_0, \epsilon} \left[\frac{1}{2\|\Sigma_\theta\|_2^2} \left\| \frac{1}{\sqrt{\alpha_t}}\left(x_t - \frac{1 - \alpha_t}{\sqrt{1 - \alpha_t}}\epsilon_t\right) - \frac{1}{\sqrt{\alpha_t}}\left(x_t - \frac{1 - \alpha_t}{\sqrt{1 - \alpha_t}}\epsilon_\theta(x_t, t)\right) \right\|^2 \right] \\ &= \mathbb{E}_{x_0, \epsilon} \left[\frac{(1 - \alpha_t)^2}{2\alpha_t(1 - \alpha_t)\|\Sigma_\theta\|_2^2} \|\epsilon_t - \epsilon_\theta(x_t, t)\|^2 \right] \\ &= \mathbb{E}_{x_0, \epsilon} \left[\frac{(1 - \alpha_t)^2}{2\alpha_t(1 - \alpha_t)\|\Sigma_\theta\|_2^2} \|\epsilon_t - \epsilon_\theta(\sqrt{\alpha_t}x_0 + \sqrt{1 - \alpha_t}\epsilon_t, t)\|^2 \right] \end{aligned} \quad (10)$$

Ho et al. (2020) propose to reweight various terms in L_{VLB} for better sample quality, to compute this objective, we generate samples $X_t \sim q(X_t|X_0)$, then train a model ϵ_θ to predict the added noise using a standard mean-squared error loss:

$$L_{simple} = \mathbb{E}_{t \sim [1, T], X_0 \sim q(X_0), \epsilon \sim N(0, I)} [\|\epsilon - \epsilon_\theta(\sqrt{\alpha_t}x_0 + \sqrt{1 - \alpha_t}\epsilon_t, t)\|^2]. \quad (11)$$

2.3 Conditional diffusion

So far, we have focused on modeling the data distribution $p(x)$. However, we are often also interested in the conditional distribution of $P(X_t|y)$, as it enables us to better investigate how different conditional information influences the generation of variable X . Begin with the score-based formulation of a diffusion model, the goal is to learn $\nabla \log P(X_t|y)$, by Bayes rules, we can get the equivalent:

$$\nabla \log P(X_t|y) = \nabla \log \left(\frac{P(y|X_t)P(X_t)}{P(y)} \right) \quad (12)$$

$$= \nabla \log P(X_t) + \nabla \log P(y|X_t) - \nabla \log P(y) \quad (13)$$

$$= \underbrace{\nabla \log P(X_t)}_{\text{unconditional score}} + \underbrace{\nabla \log P(y|X_t)}_{\text{conditional score}} \quad (14)$$

To better control the conditional information, a hyperparameter γ is introduced to

scale the gradient of the conditioning information. The score function can then be summarized as:

$$\nabla \log P(X_t|y) = \nabla \log P(X_t) + \gamma \nabla \log P(y|X_t). \quad (15)$$

Intuitively speaking, the $\gamma = 0$ the diffusion model can ignore the conditional information entirely, while a large γ value would cause the model to heavily incorporate the conditional information during sampling. In order to implement effective control over the conditional information, we use classifier-free guidance (Ho & Salimans, 2021). To get the score function under Classifier-Free Guidance, we can rearrange:

$$\nabla \log P(y|X_t) = \nabla \log P(X_t|y) - \nabla \log P(X_t). \quad (16)$$

Substituting equation (16) into equation (15) then we get:

$$\nabla \log P(X_t|y) = \nabla \log P(X_t) + \gamma(\nabla \log P(X_t|y) - \nabla \log P(X_t)). \quad (17)$$

$$= \underbrace{(1 - \gamma)\nabla \log P(X_t)}_{\text{unconditional score}} + \underbrace{\gamma\nabla \log P(X_t|y)}_{\text{conditional score}} \quad (18)$$

From Tweedie’s formula and equation 5, we can get,

$$\nabla \log p(x_t) = -\frac{1}{\sqrt{1-\bar{\alpha}_t}}\epsilon \quad (19)$$

The equation 19 means that estimating ϵ is equivalent to estimating a scaled version of the score function. So, in this paper, we model the conditional distribution of precipitation frames in the future given the past precipitation frames $\mathbf{P} = [p_1, p_2, \dots, p_M]$, we learn two sets of neural networks, $\epsilon_\theta(X_t, t)$ and $\epsilon_\theta(X_t, t, P)$, to approximate the unconditional and conditional score functions $\nabla \log P(X_t)$ and $\nabla \log P(X_t|y)$, our conditional diffusion loss function is:

$$L_{\text{condition}} = \mathbb{E}_{t \sim [1, T], X_0 \sim q(X_0), \epsilon \sim N(0, I)} [\|\epsilon - \epsilon_\theta(X_t, t, P)\|^2]. \quad (19)$$

2 Details of baseline model

2.1 Generative models of radar

DGMR holds the current state of the art in precipitation nowcasting, the generator is built with convolutional and convolutional GRU layers and it was trained with two adversarial loss and a regularization loss. The first loss is defined by a spatial discriminator, which ensures spatial consistency. The second loss is defined by a temporal discriminator, which is a 3D convolutional neural network that aims to impose temporal consistency. The regularization term encourages the prediction’s mean precipitation fields to match the mean of past precipitation amount.

We utilized Google-Colab to load the saved DGMR model and pconducted inference on our test dataset, see <https://github.com/deepmind/deepmind-research/tree/master/nowcasting>. DGMR exhibits the capability to generate forecasts up to 90 minutes. However, for the purpose of comparison, we only evaluated its performance using the first 30 minutes of forecasted results, calculating relevant metrics.

2.2 U-Net

We use a U-Net encoder–decoder model as baseline similarly to how it was used in related studies (Ayzel et al., 2020). This type of model first employs an encoder that reduces the spatial resolution using pooling and convolutional layers, while the decoder then increases the resolution by applying up-sampling and convolutional layers to the

learned patterns. To prevent gradient vanishing and share the low-level patterns of the precipitation fields, skip connections are used from the encoder to the decoder (Srivastava et al., 2015). In this paper, U-Net serves as the baseline for deterministic forecasting using deep learning.

2.3 PySTEPS

PySTEPS is an open-source Python library designed for radar precipitation forecasting and analysis, it is available at <https://github.com/pySTEPS/pysteps>. It offers a comprehensive range of algorithms, among which STEPS is a widely used precipitation nowcasting system based on ensembles, considered to be state-of-the-art of non-ML-based method. In this study, we adopt PySTEPS as a non-machine learning baseline.

3 Details of metrics

we use the M to denote number of the ensemble members, and f_m to denote the ensemble member, so the ensemble mean can be written as,

$$\bar{f} = \frac{1}{M} \sum_{m=1}^M f_m \quad (20)$$

3.1 MAE

The (spatial) mean-absolute-error (MAE) at forecast time step t between ensemble means \bar{f} and observation f_{obs} is defined as,

$$MAE_t(\bar{f}, f_{obs}) = \frac{1}{p} \sum_{p=1}^p |\bar{f} - f_{obs}| \quad (21)$$

where p indexes all the geospatial locations. And we can consider extreme value prediction accuracy under different precipitation intensities, we use an intensity mask $[f_{obs} > 4]$ and $[f_{obs} > 8]$ to get the masked prediction and observation \bar{f}_m, f_{m_obs}

$$MAE_{t,mid}(\bar{f}_m, f_{m_obs}) = \frac{1}{p} \sum_{p=1}^p |\bar{f}_m - f_{m_obs}| \quad (22)$$

3.2 Correlation

The spital correlation between ensemble mean and observation is defined as.

$$Corr_t(\bar{f}, f_{obs}) = \frac{\sum_p (\bar{f}_p - \bar{\bar{f}}_p)(f_{obs,p} - \bar{\bar{f}}_{obs,p})}{\sqrt{\sum_p (\bar{f}_p - \bar{\bar{f}}_p)^2} \sqrt{\sum_p (f_{obs,p} - \bar{\bar{f}}_{obs,p})^2}} \quad (23)$$

where $\bar{\bar{f}}_p$ means to average in space. In deployment, we flatten the prediction and observation then use the *corrcoef* function from the *NumPy* library.

3.3 Critical Success Index

The Critical Success Index (CSI) is a statistical measure that quantifies the accuracy of spatial prediction by evaluating the correct identification of specific events or outcomes.

The CSI is defined as the ratio of true positives (TP) to the sum of true positives, false positives (FP), and false negatives (FN). Mathematically, it is expressed as,

$$CSI = \frac{TP}{TP+FP+FN} \quad (24)$$

- TP represents the number of true positive outcomes, which signifies the accurate prediction of events or occurrences.
- FP corresponds to false positives, indicating instances where the event was predicted, but did not materialize.
- FN denotes false negatives, signifying cases where the event occurred but was not correctly predicted.

The CSI values range between 0 and 1, where a CSI of 1 indicates perfect spatial accuracy in prediction, implying that all positive outcomes were correctly forecasted without any false alarms. Conversely, a CSI of 0 suggests that none of the events were accurately predicted.

3.4 Continuous Ranked Probability Score

CRPS is used to evaluate the calibration and sharpness. It quantifies the discrepancy between the forecasted cumulative distribution function (CDF) and the observed CDF, defined as,

$$CRPS = \int_{-\infty}^{+\infty} [F(f_m) - 1(t \leq z)]^2 dz \quad (25)$$

where F denotes the CDF of the prediction distribution and $1(t \leq z)$ is an indicator function that is 1 if $t \leq z$ and 0 otherwise. In the case of a deterministic forecast (like Unet) the CRPS reduces to the mean absolute error (MAE).

3.5 Spread-skill ratio

The SSR evaluates the reliability of the ensemble. It is a ratio that quantifies the balance between calibration and sharpness, providing insights into the trade-off between these two critical aspects of predictive modeling.

$$SSR = \frac{Spread}{RMSE} \quad (26)$$

where the spread is defined as,

$$Spread = \sqrt{\frac{1}{P} \sum_{p=1}^P Var(f_{m,p})} \quad (27)$$

and the RMSE is defined as,

$$RMSE = \sqrt{\frac{1}{P} \sum_{p=1}^P (\bar{f} - f_{obser})^2} \quad (28)$$

4 Additional results

4.1 Skill evaluation

Figure S1 includes PySTEPS metrics calculated over the entire test dataset. Due to UNet's blurred predictions, it falls short of PySTEPS in terms of CSI8.

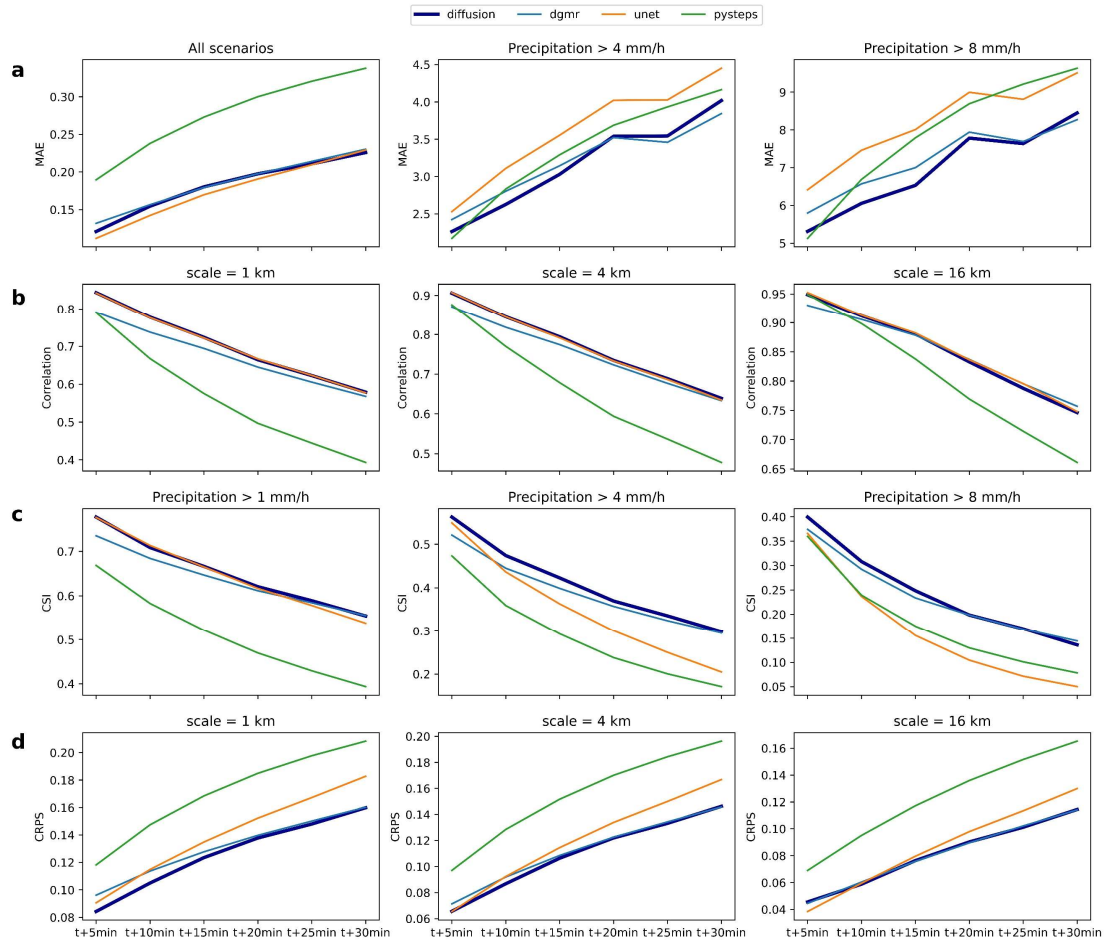


Figure S1. The full version of Figure 2, incorporating metrics for PySTEPS.

4.2 Additional case

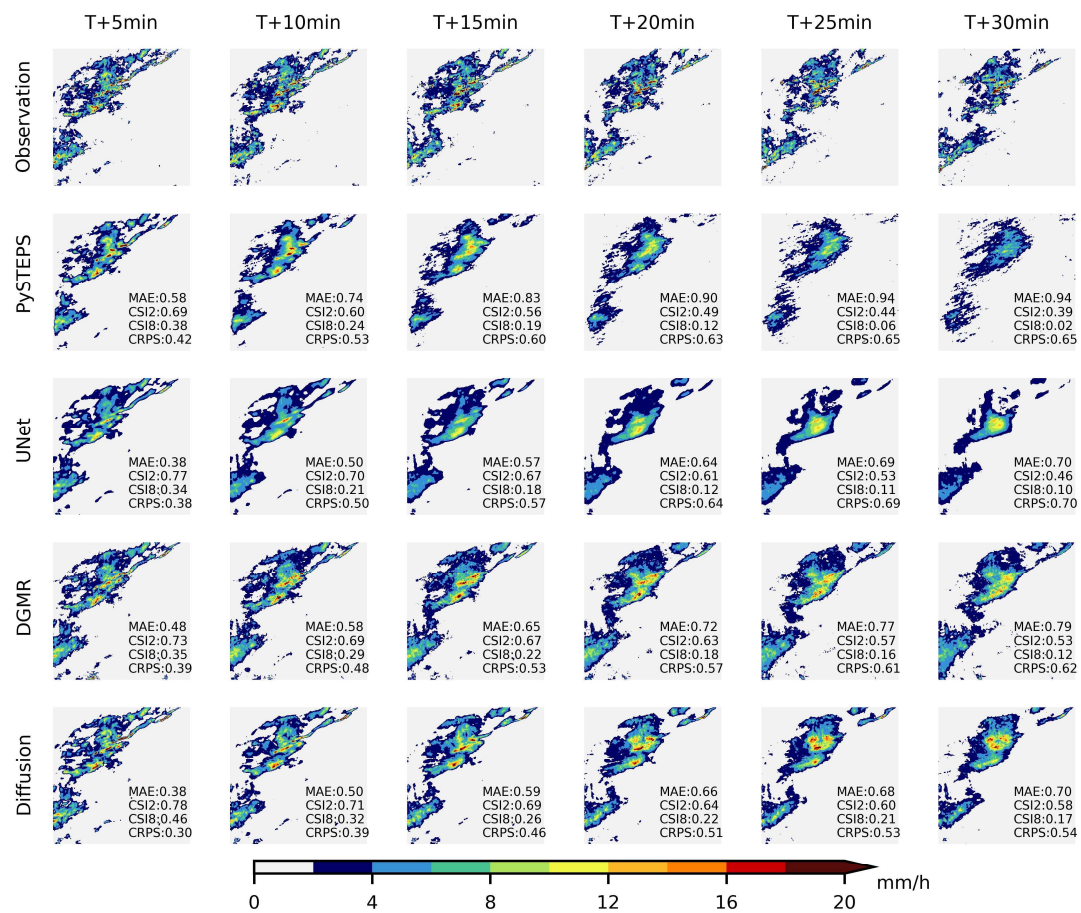


Figure S2. An additional case in Section 5.1

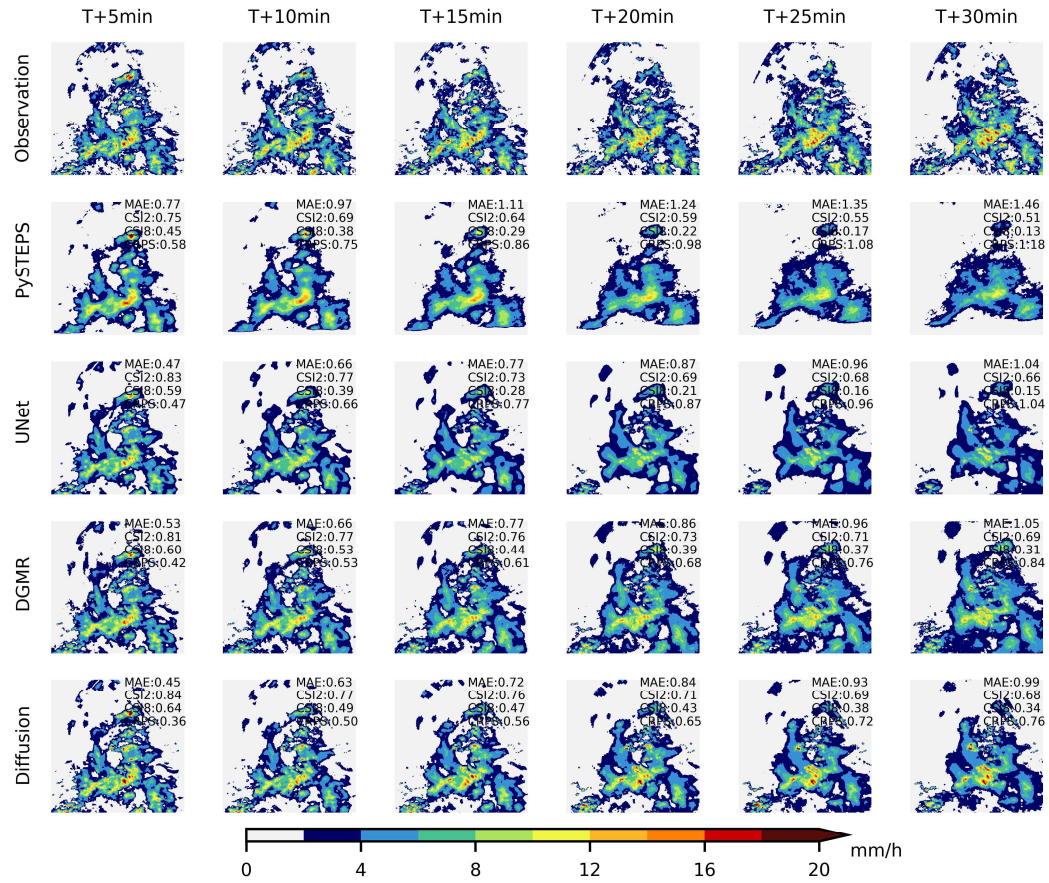


Figure S3. An additional case in Section 5.1

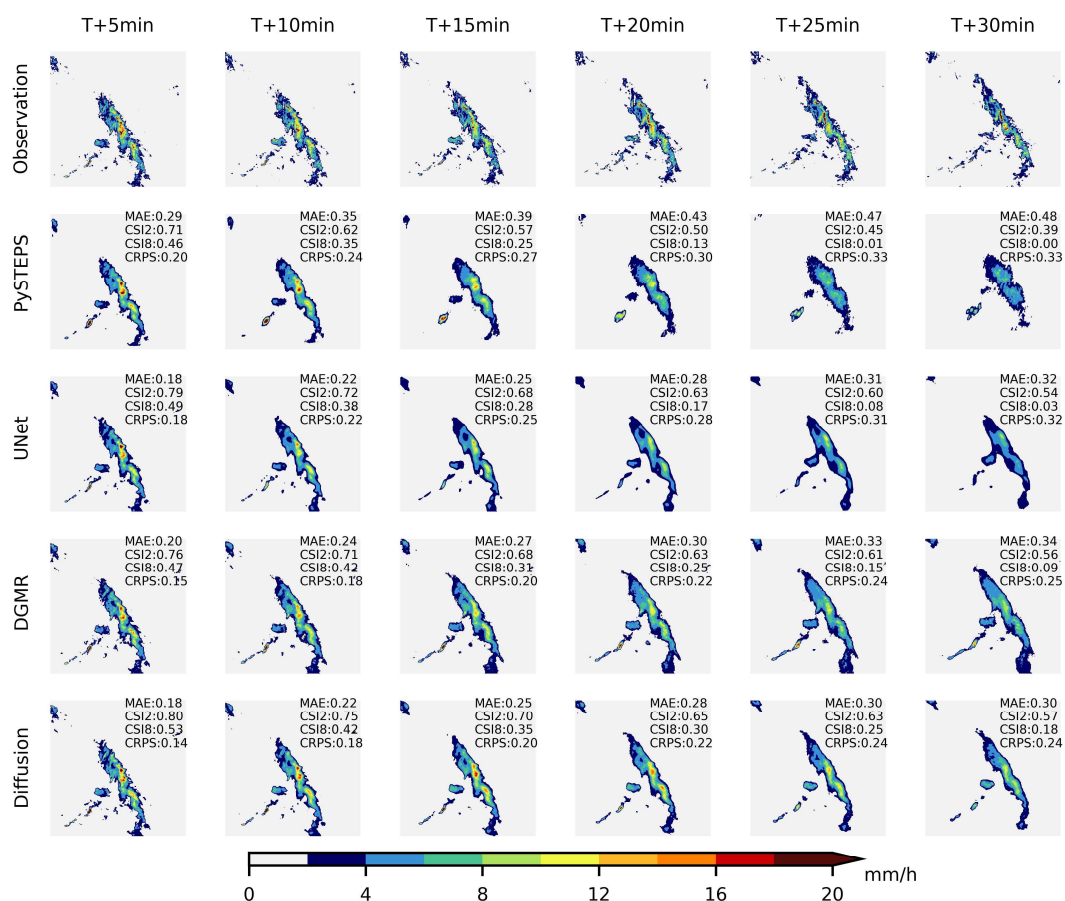


Figure S4. An additional case in Section 5.1

4.3 Reliability cases

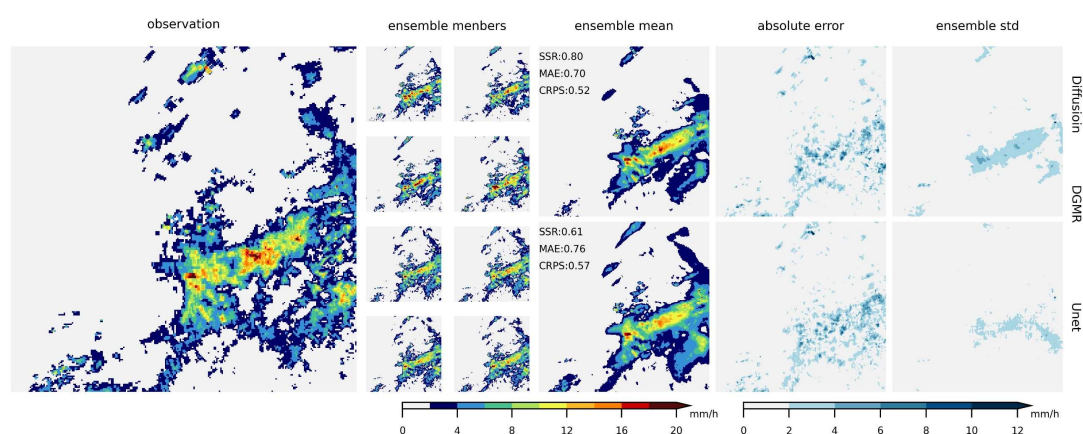


Figure S5. An additional case in Section 5.3

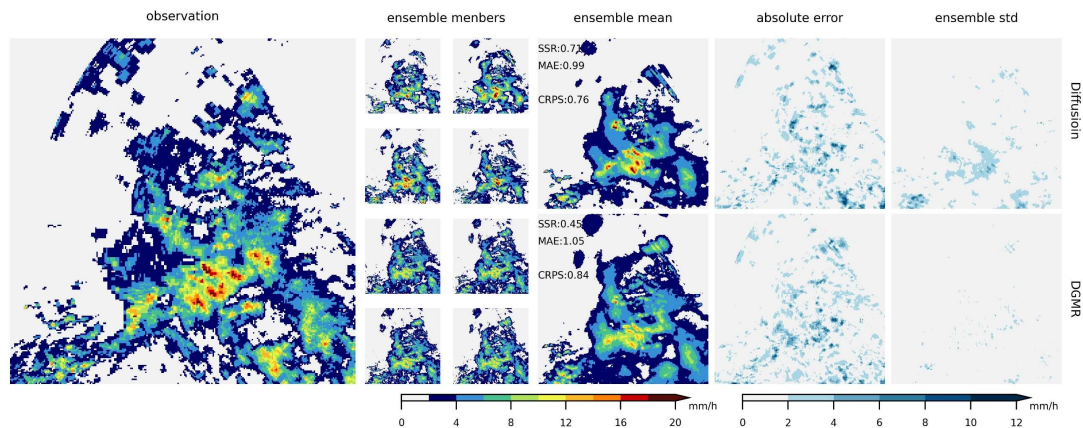


Figure S6. An additional case in Section 5.3

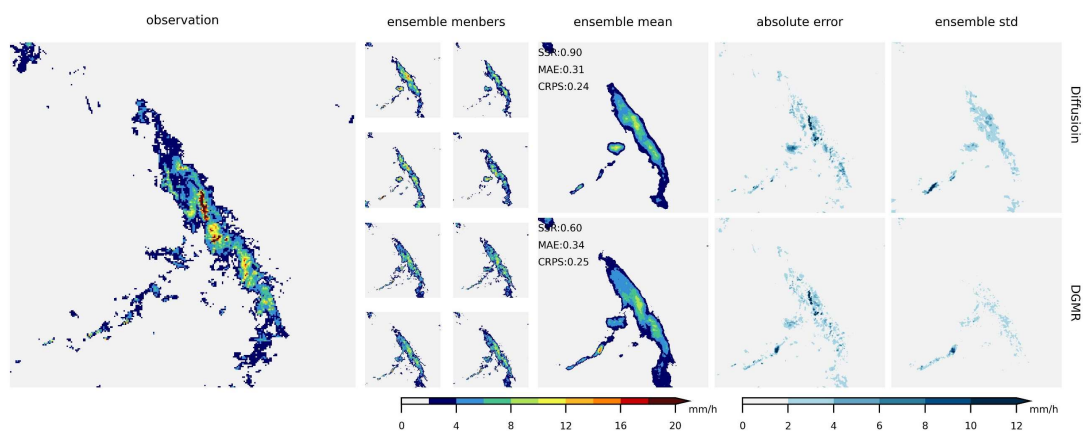


Figure S7. An additional case in Section 5.3

Reference

- Vikram Voleti, Alexia Jolicoeur-Martineau and Christopher Pal (2022). MCVD: Masked Conditional Video Diffusion for Prediction, Generation, and Interpolation. <https://doi.org/10.48550/arXiv.2205.09853>
- Understanding diffusion models: A unified perspective." arXiv preprint arXiv:2208.11970 (2022). <https://doi.org/10.48550/arXiv.2208.11970>
- Yang, L., Zhang, Z., Song, Y., Hong, S., Xu, R., Zhao, Y., ... & Yang, M. H. (2022). Diffusion models: A comprehensive survey of methods and applications. ACM Computing Surveys. <https://doi.org/10.48550/arXiv.2209.00796>