Probabilistic diffusion model for stochastic parameterization – a case example of numerical precipitation estimation

Baoxiang Pan¹, Le-Yi Wang², Feng Zhang³, Qingyun Duan⁴, Xin Li⁵, Xiaoduo Pan⁶, Xi Chen¹, Fenghua Ling⁷, Shuguang Wang⁸, Ming Pan⁹, and Ziniu Xiao¹

¹Institute of Atmospheric Physics, Chinese Academy of Sciences
²Nanjing University
³Fudan University
⁴Hohai University
⁵Institute of Tibetan Plateau Research, Chinese Academy of Sciences
⁶Institute of Tibetan Plateau Research Chinese Academy of Sciences
⁷Nanjing University of Information Science and Technology
⁸School of Atmospheric Sciences, Nanjing University
⁹University of California San Diego

December 3, 2023

Abstract

Estimating the unresolved geophysical processes from resolved geophysical fluid dynamics is the key for improving numerical weather-climate predictions. While data-driven parameterization for unresolved geophysical processes shows potential, most practices fail to capture the diversity of unresolved geophysical processes that agree with resolved geophysical fluid state. This pitfall undermines the likelihood or severity of simulated weather extremes, and erodes the fidelity of climate projections. We propose the criteria of READS (Realism, Efficiency, Adaptability, Diversity, Sharpness) for generative models to yield reasonable stochastic parameterization. We introduce probabilistic diffusion model, a non-equilibrium thermodynamics inspired deep generative modeling approach, to better meet these criteria. Using a case example of numerical precipitation estimation, we demonstrate the advantage of the proposed methodology in quickly delivering diverse and faithful estimates for the target unresolved process, as compared to other popular data-driven deterministic and stochastic methods (UNet, variational autoencoder, generative adversarial net), as well as dynamical downscaling method (WRF). We conclude that generative models, in particular, probabilistic diffusion model, can significantly enhance the representation of unresolved geophysical processes in numerical weather-climate predictions.

Probabilistic diffusion model for stochastic 1 parameterization 2 – a case example of numerical precipitation estimation 3 Baoxiang Pan¹, Leyi Wang², Feng Zhang³, Qingyun Duan⁴, Xin Li⁵, Xiaoduo 4 Pan⁵, Xi Chen¹, Fenghua Ling⁶, Shuguang Wang⁷, Ming Pan⁸, Ziniu Xiao¹ 5 ¹Institute of Atmospheric Physics, Chinese Academy of Sciences 6 ²Chongqing Institute of Big Data, Peking University ³Department of Atmospheric and Oceanic Sciences, Fudan University ⁴National Key Laboratory of Water Disaster Prevention, Hohai University ⁵Institute of Tibetan Plateau Research, Chinese Academy of Sciences ⁶Institute for Climate and Application Research, Nanjing University of Information Science and 8 9 10 11 Technology 12 ⁷School of Atmospheric Sciences, Nanjing University ⁸Scripps Institution of Oceanography, University of California San Diego 13 14

To learn stochastic parameterization, one should steer generative models toward the requirements of ensemble forecasts. We propose criteria of READS (Realism, Efficiency, Adaptability, Diversity, Sharpness) for data-driven stochastic parameterization. In case example of numerical precipitation estimation, probabilistic diffusion model well meets the READS criteria.

Key Points:

15

Corresponding author: Baoxiang Pan, panbaoxiang@lasg.iap.ac.cn

22 Abstract

Estimating the unresolved geophysical processes from resolved geophysical fluid dynam-23 ics is the key for improving numerical weather-climate predictions. While data-driven 24 parameterization for unresolved geophysical processes shows potential, most practices 25 fail to capture the diversity of unresolved geophysical processes that agree with resolved 26 geophysical fluid state. This pitfall undermines the likelihood or severity of simulated 27 weather extremes, and erodes the fidelity of climate projections. We propose the crite-28 ria of READS (Realism, Efficiency, Adaptability, Diversity, Sharpness) for generative mod-29 els to yield reasonable stochastic parameterization. We introduce probabilistic diffusion 30 model, a non-equilibrium thermodynamics inspired deep generative modeling approach, 31 to better meet these criteria. Using a case example of numerical precipitation estima-32 tion, we demonstrate the advantage of the proposed methodology in quickly delivering 33 diverse and faithful estimates for the target unresolved process, as compared to other 34 popular data-driven deterministic and stochastic methods (UNet, variational autoencoder, 35 generative adversarial net), as well as dynamical downscaling method (WRF). We con-36 clude that generative models, in particular, probabilistic diffusion model, can significantly 37 enhance the representation of unresolved geophysical processes in numerical weather-climate 38 predictions. 39

⁴⁰ Plain Language Summary

⁴¹ "Life is a gorgeous robe, crawling with lice", so said Eileen Chang, a Chinese writer
⁴² who enjoyed depicting the awkward discrepancies between ideal and reality. Same metaphor
⁴³ applies to climate models, rooted in physical principles of fluid dynamics and thermo⁴⁴ dynamics, rife with empirics making up the missing components. We use generative AI
⁴⁵ to make up the missing components in climate models, achieving realistic and informa⁴⁶ tive simulations of unresolved climate processes, i.e., precipitation.

47 **1** Introduction

Geophysical fluid dynamics operates across a continuous spectrum of spatiotemporal scales, ranging from micro-scale turbulences to synoptic-scale planetary waves. Their numerical solvers, coming with finite resolution, set a distinction between resolved dynamics and unresolved physical processes, with the latter being approximated as empirical functions of the former. This approximation, known as parameterization, is the source of error in numerical weather and climate predictions (Stensrud, 2009).

Typically, parameterization schemes are deterministic functions, providing a unique 54 tendency accounting for the grid-scale impact of subgrid physical processes in numer-55 ical modeling of geophysical fluid dynamics. However, as we do not explicitly resolve the 56 subgrid physical processes, a probabilistic formulation is advocated (Berner et al., 2017; 57 T. Palmer, 2019): the impact of subgrid physical processes should be described by a prob-58 ability distribution function conditioning on the resolved geophysical fluid dynamics. This 59 probabilistic formulation enables a rigorous and consistent characterization of unresolved 60 physical processes across model resolutions (Sakradzija et al., 2016). Also, it allows subgrid-61 scale noise to trigger crucial circulation regime transitions, supporting reliable proba-62 bilistic forecasts (T. N. Palmer et al., 2009). 63

To make accurate probabilistic representation of unresolved physical processes in numerical weather-climate models, existing efforts proceed along the following three lines. The first, which is straightforward yet lacks theoretical warrant, is to pre-define a perturbation to the parameters, functions, or outputs of deterministic parameterization schemes (Dorrestijn et al., 2013). The second, which is solid yet costly and restricted in scope, is to compute statistics of the equilibrium states of the considered process, via statistical mechanics analysis (Plant & Craig, 2008). The third, which promises remarkable accuracy, efficiency, and flexibility, is to approximate the probability distribution of unresolved physical processes by *learning* from high fidelity data, such as high-resolution simulations or observations (Gagne et al., 2020; Ravuri et al., 2021; Harris et al., 2022).

We proceed along the third line as motivated by the recent breakthrough of prob-74 abilistic machine learning, in particular, probabilistic diffusion models (Sohl-Dickstein 75 et al., 2015; Ho et al., 2020; Song et al., 2020). Probabilistic diffusion models learn to 76 approximate probability distributions in an iterative manner, achieving unprecedented 77 fitting capacity and controlling flexibility in generative modeling tasks. Using a case ex-78 79 ample of numerical precipitation estimation, we specify five requirements for developing data-driven stochastic parameterization schemes. We develop Diffusion based Pre-80 cipitation estimator, dubbed DiP, and demonstrate its unique advantages in meeting these 81 requirements, as compared to existing data-driven deterministic and stochastic param-82 eterization schemes, as well as high resolution dynamical simulation method. 83

⁸⁴ 2 Problem setup and model requirements

We consider a case example of numerical precipitation estimation: given geophys-85 ical fluid dynamics resolved to a finite spatiotemporal resolution, the goal is to estimate 86 the accompanying precipitation process. The challenge lies in that, precipitation results 87 from a complicated chain of processes that are mostly unresolved in numerical models 88 (Tapiador et al., 2019). Any error along this simulation chain may distort the location, 89 timing, or quantity of the precipitation estimate, rendering the estimate useless, even mis-90 leading (Pan, Hsu, AghaKouchak, & Sorooshian, 2019; Pan, Hsu, AghaKouchak, Sorooshian, 91 & Higgins, 2019; Chen & Wang, 2022). 92

Here we consider the region of East and Southeast Asia $(0^{\circ}-40^{\circ}N, 100^{\circ}E-140^{\circ}E)$, 93 where precipitation is driven by diverse circulation regimes. We use the following resolv-94 able dynamical variables to infer precipitation: key primitive variables (meridional and 95 zonal wind velocity, temperature, specific humidity, and geopotential height) at 3 pres-96 sure levels (1000/850/500 hPa), and crucial surface level variables (sea level pressure, 97 surface pressure, surface temperature, and total column precipitable water). These data 98 are obtained by blending observations with short-range weather forecasts to faithfully 99 represent historical circulation states. The data are from the Climate Forecast System 100 Reanalysis project (Saha et al., 2006), coming at spatiotemporal resolution of $0.5^{\circ}/1$ hour 101 for Year 1979-2022. Besides these dynamical variables, we also consider 0.1° elevation 102 data as extra, static predictor. These dynamical and static predictor variables are to-103 gether denoted as \mathbf{x} . Precipitation, as our predict and variable, is denoted by \mathbf{y} . The data 104 are from the Multi-Source Weighted-Ensemble Precipitation product (Beck et al., 2019), 105 which merges gauge, satellite, and reanalysis precipitation records to achieve optimal qual-106 ity. The data come at a $0.1^{\circ}/3$ -hourly resolution for same period. 107

Our objective is to approximate the conditional distribution of $p(\mathbf{y}|\mathbf{x})$, based on favorably large amount of $\{\mathbf{x}, \mathbf{y}\}$ paired data samples. This problem setup differs from a deterministic regression problem setup, which has been widely adopted for learning parameterization schemes (Yu et al., 2023; Wang & Tan, 2023). Specifically, in both deterministic and probabilistic formulations, we design a learning machine Θ , for which the optimal parameter θ^* is obtained by maximizing the overall likelihood of the reference data:

$$\theta^* = \operatorname*{argmax}_{\theta} \sum_{i} \log p_{\theta}(\mathbf{y}_i | \mathbf{x}_i) \tag{1}$$

In a deterministic regression problem setup, given any \mathbf{x} , the learning machine yields the most plausible \mathbf{y} : $\hat{\mathbf{y}} = \underset{\mathbf{y}}{\operatorname{argmax}} \log p_{\theta^*}(\mathbf{y}|\mathbf{x})$. This is often achieved by pre-assuming the distributional form of p_{θ} . For instance, assuming $p_{\theta}(\mathbf{y}|\mathbf{x}) := \mathcal{N}(\mu_{\theta}(\mathbf{x}), \sigma_{\theta}^{2}(\mathbf{x}))$, we have $\hat{\mathbf{y}} = \mu_{\theta^{*}}(\mathbf{x})$, where $\theta^{*} = \underset{\theta}{\operatorname{argmin}} \sum_{i} \left[\frac{(\mu_{\theta}(\mathbf{x}_{i}) - \mathbf{y}_{i})^{2}}{\sigma_{\theta}^{2}(\mathbf{x}_{i})} + \log \sigma_{\theta}^{2}(\mathbf{x}_{i}) \right]$. Such an assumption offers huge computation convenience, yet, it comes with two deficits. First, a predefined distributional form often poorly fits a richly structured target physical process. Second, a deterministic formulation precludes interaction between subgrid noise and resolved dynamics, resulting in biased weather-climate predictions (Hardiman et al., 2022).

In a probabilistic modeling setup, given any \mathbf{x} , the learning machine outputs plau-123 sible y samples: $\hat{\mathbf{y}} \sim p_{\theta^*}(\mathbf{y}|\mathbf{x})$, where p_{θ^*} is a learned distribution subject to no pre-124 defined probability distribution form. This probabilistic formulation allows us to bypass 125 the two deficits of a deterministic formulation, yet, it comes with its own challenges and 126 requirements. To realize the potential, one must steer the learning machine toward ver-127 ifiable goals of stochastic parameterization, which are quantified in ensemble forecast prac-128 tices. We hence suggest the following five criteria for p_{θ^*} based on the requirements of 129 ensemble forecast: 130

131	• Realism : samples from the estimated conditional probability distribution should
132	be indistinguishable from observational samples, regarding either their structure
133	or functionality. This requirement ensures that an accurate probability value can
134	be assigned to the realized observations, either for training or evaluation purposes.
135	Also, the generated samples can fit into the geophysical modeling pipeline, and
136	be as useful as observations for a wide range of subsequent tasks.
137	• Efficiency: a solid approach for developing parameterization schemes is to con-
138	sider each of the possible ways that the subgrid scale process evolve under the grid-
139	scale constraint, to compute the probability of each such "configuration" in the
140	equilibrium ensemble, and generate samples accordingly. This requires excessive
141	human effort and computational resources. Here, we expect a well-trained p_{θ^*} to
142	efficiently generate multiple samples of plausible subgrid physical processes, at least
143	several orders faster than directly resolving the subgrid scale process.
144	• Adaptability: the interaction of subgrid scale physics and large scale dynamics
145	often results in organized weather schemes across scales, ranging from local con-
146	vection to weather fronts. Correspondingly, the model is preferred to automati-
147	cally identify and apply to these organized weather schemes, rather than work-
148	ing at individual computing grids or fixed computing time steps.
149	• Diversity: the estimated conditional probability distribution should cover all plau-
150	sible outcomes, rather than a limited subset of modes. This ensures that all ob-
151	served states are within the <i>cone</i> of model simulations, particularly for extremes.
152	• Sharpness: the estimated conditional probability distribution should generate
153	samples that are faithful to the conditioning information, maximizing the sharp-
154	ness of the simulated distribution. Note that this requirement naturally confronts
155	the diversity requirement: an overly constrained probability estimate may fail to
156	encapsulate observations, which is unreliable; an overly dispersed probability es-
157	timate may lack clear distinction from climatology, which is uninformative. We
158	must carefully balance sharpness and diversity, so that the probability estimate
159	faithfully reflects the intrinsic stochasticity of the considered process.

We coin the term READS by concatenating the initial letters of the five criteria above. Below we introduce DiP, Diffusion based Precipitation estimator, and demonstrate its unique advantage in meeting the READS requirements, as compared to existing deterministic/probabilistic data-driven approaches, as well as high-resolution dynamical simulation approach.

¹⁶⁵ 3 Diffusion based Precipitation estimator (DiP)

3.1 A primer on probabilistic machine learning

The method we develop here falls into the scope of probabilistic machine learning, 167 which applies probability theory to design learning machines that make predictions as 168 probability distributions. Since the target distribution we try to approximate adheres 169 to no pre-defined closed form, a common strategy is to learn a mapping between the tar-170 get distribution and a tractable latent distribution, i.e., standard Gaussian. After learn-171 ing the mapping from optimally large amount of data, we can pass samples from the la-172 tent distribution through the trained model to obtain target distribution samples, hence 173 inferring this target distribution. The key challenge is that, we lack point-to-point cor-174 respondences between samples from the target distribution and the latent distribution, 175 hence lacking straightforward supervision signals to enable learning (Ruthotto & Haber, 176 2021). A popular solution is to build bijective mapping between the target distribution 177 and the latent distribution, therefore establishing correspondences. Probabilistic diffu-178 sion models excel in this task by establishing the bijection in an iterative manner. Be-179 low we outline how this is achieved. Mathematical and implementation details are given 180 in Supporting Information S1. 181

¹⁸² 3.2 Basics

166

¹⁸³ Diffusion model approximates a target distribution $p(\mathbf{y})$ by reversing a Gaussian ¹⁸⁴ process (Fig. 1): the forward Gaussian process turns $p(\mathbf{y})$ into standard Gaussian $\mathcal{N}(\mathbf{0}, \mathbf{I})$ ¹⁸⁵ (Fig. 1a, Eq. 2); we learn to iteratively reverse this Gaussian process, mapping $\mathcal{N}(\mathbf{0}, \mathbf{I})$ ¹⁸⁶ to $p(\mathbf{y})$ (Fig. 1b-d), hence achieving generative modeling. Following D. Kingma et al. (2021), ¹⁸⁷ the forward Gaussian process is pre-defined as:

$$p(\mathbf{z}_t | \mathbf{y}) := \mathcal{N}(\alpha_t \mathbf{y}, \sigma_t^2 \mathbf{I})$$
(2)

Here \mathbf{z}_t is latent variable indexed by $t \in [0, 1]$, α_t / σ_t is monotonically decreasing/increasing function of t, strictly bounded by [0, 1]. Eq. 2 therefore bridges $p(\mathbf{y}) = p(\mathbf{z}_0|\mathbf{y})$ and $\mathcal{N}(\mathbf{0}, \mathbf{I}) = p(\mathbf{z}_1|\mathbf{y})$ (Fig. 1a). We reverse Eq. 2 to turn $\mathcal{N}(\mathbf{0}, \mathbf{I})$ into $p(\mathbf{y})$, using a chain of variational distributions (Fig. 1b):

$$p_{\theta}(\mathbf{z}_{t_{i-1}}|\mathbf{z}_{t_i}) = \mathcal{N}(\mu_{\theta}(\mathbf{z}_{t_i}), \mathbf{\Sigma}_{\theta}(\mathbf{z}_{t_i})), \quad i \in [1, T]$$
(3)

Here $0 = t_0 < t_1 < t_2 < ... < t_T = 1$ is arbitrary discretization of time; $\{\mu_{\theta}, \Sigma_{\theta}\}$ are

learnable mean vector and covariance matrix, trained by maximizing the overall data like lihood (Supporting Information S1.1):

$$\log p_{\theta}(\mathbf{y}) = \mathbb{E}_{p(\mathbf{z}_{0}|\mathbf{y})} \log p(\mathbf{y}|\mathbf{z}_{0}) - D_{\mathrm{KL}} \left(p(\mathbf{z}_{1}|\mathbf{y}) || p(\mathbf{z}_{1}) \right) - \sum_{i=1}^{T} \mathbb{E}_{p(\mathbf{z}_{t_{i}}|\mathbf{y})} D_{\mathrm{KL}} \left(p(\mathbf{z}_{t_{i-1}}|\mathbf{z}_{t_{i}},\mathbf{y}) || p_{\theta}(\mathbf{z}_{t_{i-1}}|\mathbf{z}_{t_{i}}) \right)$$
(4)

Given Eq. 2, to maximize Eq. 4 is approximately equivalent to minimizing the Fisher divergence between the data and model distributions (Supporting Information S1.2):

$$\theta^* = \operatorname*{argmax}_{\theta} \log p_{\theta}(\mathbf{y}) \approx \operatorname*{argmin}_{\theta} \sum_{i=1}^{T} \mathbb{E}_{p(\mathbf{z}_{t_i}|\mathbf{y})} \left\| \nabla \log p(\mathbf{z}_{t_i}|\mathbf{y}) - \epsilon_{\mathrm{NN}_{\theta}}(\mathbf{z}_{t_i}) \right\|_2$$
(5)

Here $\epsilon_{NN_{\theta}}$ is a neural network parameterization of $\nabla \log p(\mathbf{z}_{t_i}|\mathbf{y})$, known as the score func-

tion. Based on the learned score estimates, we can derive $p_{\theta}(\mathbf{z}_{t_{i-1}}|\mathbf{z}_{t_i}) = \mathcal{N}(\mu_{\theta}(\mathbf{z}_{t_i}), \boldsymbol{\Sigma}_{\theta}(\mathbf{z}_{t_i}))$ (Supporting Information S1.2) and sample it, starting with $p(\mathbf{z}_1) = \mathcal{N}(\mathbf{0}, \mathbf{I})$, ending with

 $n(\mathbf{z}_0) \approx n(\mathbf{y})$

198 $p(\mathbf{z}_0) \approx p(\mathbf{y}).$



Figure 1. Overview of diffusion model. We map target distribution (synoptic-scale precipitation field, **a** left) to a same dimensional standard Gaussian distribution (**a** right) through a pre-defined Gaussian process (**a** bottom, Eq. 2). Color denotes probability distribution function value for an individual precipitation field pixel (here we select the center pixel) through diffusion time t = [0, 1], lines show the diffusion trajectories of individual pixels for randomly selected samples, matrix plots show the noisified precipitation field (sample of Typhoon Lekima, 0000 UTC 09 August 2019, centered at 26.5°N, 114.4°E) across diffusion time (**a** top). We approximate the target distribution by reversing the Gaussian process, using a series of variational distributions (**b**, Eq. 3), which are trained by maximizing the data likelihood (Eq. 4-5). We include conditioning information to approximate conditional distribution of a same Typhoon event (**c**). We apply *classifier-free guidance* to control the impact of the conditioning information versus the latent variable in explaining the variability of the target variable for the same event (**d**, Eq. 6). By enhancing the guidance strength ω , we suppress the variance of the resulting conditional probability distribution (**c**/**d** right). The plots are supported by logarithm transformed precipitation observational data for Year 2019, and the trained diffusion models.

¹⁹⁹ 3.3 Conditioning

To generate \mathbf{y} samples that are faithful to the conditioning information \mathbf{x} , we need 200 to approximate the conditional distribution $p(\mathbf{y}|\mathbf{x})$. To achieve this, we include \mathbf{x} dur-201 ing training and sampling (Fig. 1c-d). A direct inclusion of \mathbf{x} does not specify the im-202 pact of **x** versus \mathbf{z}_t in explaining the variability of **y** (Fig. 1c, Holmes & Walker 2017). 203 To tackle the potential misspecification, and having \mathbf{x} effectively control the learned dis-204 tribution, we resort to *classifier-free guidance* (Ho & Salimans, 2022, Fig. 1d): we learn 205 two sets of neural networks: $\epsilon_{\rm NN}(\mathbf{z}_{t_i}) / \epsilon_{\rm NN}^c(\mathbf{z}_{t_i}, \mathbf{x})$, so to approximate the unconditional/conditional 206 scores: $\nabla \log p(\mathbf{z}_{t_i}|\mathbf{y}) / \nabla \log p(\mathbf{z}_{t_i}|\mathbf{y}, \mathbf{x})$. Based upon these two sets of score estimates, 207 we compose score estimators for synthetic distributions $p_{\omega}(\mathbf{z}_{t_i}|\mathbf{x},\mathbf{y}) \propto p(\mathbf{x}|\mathbf{y},\mathbf{z}_{t_i})^{\omega} p(\mathbf{z}_{t_i}|\mathbf{y})$: 208

$$\nabla \log p_{\omega}(\mathbf{z}_{t_i}|\mathbf{y}, \mathbf{x}) = \nabla \log p(\mathbf{z}_{t_i}|\mathbf{y}) + \omega \nabla \log p(\mathbf{x}|\mathbf{y}, \mathbf{z}_{t_i})$$

= $\nabla \log p(\mathbf{z}_{t_i}|\mathbf{y}) + \omega (\nabla \log p(\mathbf{z}_{t_i}|\mathbf{x}, \mathbf{y}) - \nabla \log p(\mathbf{z}_{t_i}|\mathbf{y}))$ (6)
 $\approx \epsilon_{\mathrm{NN}}(\mathbf{z}_{t_i}) + \omega (\epsilon_{\mathrm{NN}}^c(\mathbf{z}_{t_i}, \mathbf{x}) - \epsilon_{\mathrm{NN}}(\mathbf{z}_{t_i}))$

Here ω is guidance scale coefficient, balancing the diversity and sharpness of the learned conditional distribution:

 $\omega = 1$: assuming impact of **x** has been perfectly accounted by $\epsilon_{\text{NN}}^c(\mathbf{z}_{t_i}, \mathbf{x})$ (Fig. 1c).

 $\omega < 1$: suppressing impact of **x**, pervading the distribution toward climatology.

 $\omega > 1$: raising impact of **x**, sharpening the distribution toward more likely values (Fig. 1d).

We now apply score estimates of $p_{\omega}(\mathbf{z}_{t_i}|\mathbf{x}, \mathbf{y})$ to sample $p(\mathbf{y}|\mathbf{x})$, following a same strategy described in Sec. 3.2. The value of ω is empirically determined based on the probabilistic forecasting skill of its resulting model.

3.4 Baselines and implementation details

We compare the DiP methodology with popular deterministic and stochastic data-

- driven methods and moderate/high resolution dynamical simulation method, including:
- **UNet**: a de-facto choice for image-to-image regression tasks, using neural network consisting symmetric convolution and deconvolution blocks (Ronneberger et al., 2015).
- Conditional variational autoencoder (CVAE): a probabilistic deep learn ing method that maximizes a lower bound of data likelihood to learn latent variable model for a target conditional distribution (D. P. Kingma & Welling, 2013;
 Pan et al., 2022).
- Conditional generative adversarial net (CGAN): a probabilistic deep learning method in which a generative network learns to approximate a target conditional distribution, under the guidance of a discriminative network that distinguishes generated samples and true samples (Goodfellow et al., 2014; Pan et al., 2021; Ravuri et al., 2021).
- CFS reanalysis precipitation product (CFSR): an optimized combination of CMAP (CPC Merged Analysis of Precipitation), daily gauge observations, and CFS background 6-hourly precipitation analysis (Saha et al., 2006).
- Dynamical downscaling using WRF: refining coarsely resolved climate pro cesses via high resolution numerical geophysical fluid dynamics solver and accom panying parameterization schemes, using Advanced Research Version 4.2 of Weather
 Research and Forecasting (WRF-ARW V4.2, Skamarock et al. 2019).
- For all the data-driven models, including DiP, we use data from 1979-2016/2017-2018/2019-2022 for training/validation/test. Considering the computation cost and the

characteristic scale of atmospheric dynamics, all the data-driven models operate at a synoptic scale $(8^{\circ} \times 8^{\circ})$: we randomly crop paired predictor and predictand field data within the study region for model training. The model structures, hyper-parameter setups, and training details are given in Supporting Information S2.

²⁴⁵ **3.5 Evaluation**

We verify models' performances using a suite of skill metrics corresponding to the 246 READS criteria. We apply Human eYe Perceptual Evaluation (HYPE, Zhou et al. 2019) 247 and power spectrum analysis to determine models' sample fidelity. We use Pearson cor-248 relation coefficient (r) and Root Mean Squared Error (RMSE) between observations and 249 models' ensemble mean estimations to quantify models' deterministic prediction skills. 250 We apply Continuous Ranked Probabilistic Skill (CRPS) to measure the accuracy of the 251 predicted probabilities and the sharpness of the forecast distribution. We compute model's 252 skill spread correlation (SSC) to quantify the reliability of a model's uncertainty esti-253 mates. We compute the ratio that observations falls into model's ensemble intervals (CR). 254 We record the computing time of the considered models. All the skill metrics are com-255 puted across spatial scales from 0.1° to 2° by aggregating neighbourhood grids. For de-256 tails, see Supporting Information S3. 257

258 4 Results

259 4.1 Case study

We start with a case example to compare models' performances. We consider the 260 storm process associated with Typhoon Lekima, which ranks as the third costliest ty-261 phoon in Chinese history. We show $8^{\circ} \times 8^{\circ}$ observed and simulated precipitation rate 262 maps along the typhoon trajectory (Fig. 2). Here, observations (Fig. 2a) present a clear 263 ring structure of intense precipitation surrounding the typhoon eye before landing (0000 264 UTC 04 August 2019 - 0000 UTC 08 August 2019), with maximum precipitation rate 265 reaching 100 mm/h. The eyewall structure gradually dissipates through two landings (1800 266 UTC 09 August and 1200 UTC 11 August), leaving a tightly curved rainband wrapping 267 into a relatively well-defined centre. 268

The large-scale patterns of precipitation estimates from the data-driven models (Fig. 2b-269 e) and CFS reanalysis (Fig. 2f) roughly agree with observations (Fig. 2a), due to a shared 270 circulation constraint from CFS reanalysis. For WRF dynamical downscaling (Fig. 2g), 271 despite careful spectral nudging, the results do not strictly follow the observed typhoon 272 trajectory, particularly after landing (1800 UTC 09 August). This is due to the chaotic 273 nature of geophysical fluid dynamics. The fine-scale structure differs significantly among 274 models: DiP (Fig. 2b) produces the most realistic small-scale details, creating a clear eve-275 wall structure and associated spiral rainband, with intense precipitation matching ob-276 servations at relatively correct locations. CGAN (Fig. 2c) can generate intense precip-277 itations surrounding the typhoon eye. Yet, the estimates come with poor spatial struc-278 ture, with neighboring grids loosely correlated, and the rainband barely depictable. CVAE 279 (Fig. 2d) and UNet (Fig. 2e) offer similar, blurry estimates, failing to distinct charac-280 teristic typhoon eyewall and rainband structures. Besides, both models miss precipita-281 tion extremes, with maximum precipitation estimates below 30 mm/h. CFS reanalysis 282 (Fig. 2f) shares similar drawbacks as CVAE and UNet, largely due to biases from the 283 assimilated data sources and errors from precipitation related model parameterization 284 schemes. WRF simulation (Fig. 2g) makes overly confined, extremely intense (approx-285 imately 150 mm/h) precipitation estimates, following the finely resolved, yet potentially 286 misaligned circulation state estimates. 287

We further inspect the probabilistic models (DiP, CGAN, and CVAE) through the lens of the READS requirements (Sec. 2). For an individual snapshot of precipitation



Figure 2. Observed and simulated $8^{\circ} \times 8^{\circ}$ precipitation rate maps along the trajectory of Typhoon Lekima, from 0000 UTC 04 August 2019 to 0000 UTC 12 August 2019. a: precipitation observations from MSWEP. b-d: randomly selected samples of ensemble precipitation estimates using DiP/CGAN/CVAE. e: deterministic precipitation estimates using UNet. f: CFS reanalysis precipitation with resolution of 0.2° . g: precipitation estimates using WRF dynamical simulation, with resolution of ~ 3 km. The typhoon trajectory from WRF simulation considerably diverges from observations after the first landing (1800 UTC 09 August). For after landing results, we show precipitation rate maps surrounding WRF simulated typhoon center.

estimate centering around 22.7°N, 125.9°E at 0000 UTC 06 August 2019, we show models' ensemble members, ensemble mean and standard deviation, ensemble mean absolute error, as well as radial/orientation averaged power spectrum (Fig. 3). We compute
a suite of skill metrics corresponding to the READS requirements.

• **Realism**: we measure human climate experts' error rate in detecting observation 294 from model estimates: for DiP/CGAN/CVAE, 3/1/0 out of 5 climate scientist eval-295 uators fail to detect the observation from 15 randomly generated model estimates, 296 suggesting the optimal spatial coherency of DiP estimates. Additionally, we in-297 spect the spatial structure of precipitation estimates by computing their average 298 spectrum power as function of spatial frequency and orientation: DiP and CGAN 299 well reproduce the spatial variability across spatial scales and orientations. Mean-300 while, WRF significantly overestimates spatial variability; CVAE, UNet and CFSR 301 significantly underestimate spatial variability for high spatial frequency and all 302 orientations. 303

- Efficiency: all the probabilistic models demonstrate advantageous efficiency compared to high-resolution numerical simulation: DiP/CGAN/CVAE generate 100 member ensemble estimates of 0.1° precipitation field within approximately 100/2/2
 seconds on a NVIDIA GeForce RTX 4090 GPU. Here, DiP is two-orders slower
 than CGAN and CVAE due to its iterative generation nature. As a comparison,
 a deterministic WRF simulation takes around 5 hours in a 32-core CPU machine.
- Adaptability: data-driven models are often reported to struggle with extremes, due to unreasonable learning objective setups, as well as approximation, optimization, and statistical errors. While the typhoon case we consider here is featured by extreme precipitation, DiP successfully reproduces the maximum precipitation rate and characteristic typhoon rainfall structures, suggesting its adaptability for extreme cases. We further report models' performances for various weather schemes in Sec. 4.2.
- **Diversity-Sharpness tradeoff**: we measure the diversity of models' ensemble 317 estimates by computing the percentage that a grid point observation falls into model's 318 ensemble interval. Here, 80.5%/53.6%/29.7% grid point observations are within 319 the 16-member ensemble interval from DiP/CGAN/CVAE. Grid points where ob-320 servations fall above/below the ensemble interval are stippled with red/black. These 321 results suggest the peculiar advantage of DiP in delivering broad range of plau-322 sible outcomes. We further investigate model's sharpness subject to a "proper" 323 level of diversity. By "proper", we mean that the probability estimate accurately 324 reflects the intrinsic stochasticity of the considered process, which is not directly 325 measurable and requires statistical inference. A good indicator is how model's en-326 semble spread aligns with model's skill. DiP achieves the highest spread-skill cor-327 relation, assigning high/low forecast uncertainty estimates to predictions with high/low 328 errors. We further consider the spatial correlation between the ensemble mean es-329 timate and observation, as well as the mean absolute error between each ensem-330 ble member and observation. The high skill values of DiP suggest that its ensem-331 ble dispersion centers around observation, requiring no ensemble pruning. Finally, 332 we report models' continuous ranked probability scores, which considers both pre-333 diction diversity and sharpness. DiP achieves the optimal performance under this 334 proper scoring rule (Gneiting & Raftery, 2007). 335
 - 4.2 Skill evaluation

336

We evaluated models' overall performances using test set data from 2019 to 2022. We report a suite of deterministic and probabilistic skill metrics for the considered models in Fig. 4.



Figure 3. Precipitation estimates centering around 22.7° N, 125.9° E at 0000 UTC 06 August 2019, using DiP (a), CGAN (b), and CVAE (c). The columns show models' ensemble members, ensemble mean, ensemble standard deviation, ensemble mean absolute error, grid points where observation is not encapsulated by ensemble spread (red/black stipple for under/over estimation, background colored based on observation), and radial/orientation averaged power spectrum for observation and all the considered models, including DiP, CGAN, CVAE, UNet, CFS reanalysis, and WRF. The following skill metrics are computed. HYPE: human climate experts' error rate in detecting observation from model estimates; r: spatial correlation between model ensemble mean estimate; CRPS: continuous ranked probabilistic score of model ensembles; SSC: spread-skill correlation, where spread is represented using ensemble standard deviation, and skill is represented using model ensemble mean absolute error; CR: coverage ratio, which represent the percentage that grid observation falls into the coverage of ensemble spread.

For deterministic evaluation, we compute the correlation coefficient (r, Fig. 4a) and 340 the root mean squared error (RMSE, Fig. 4b) between observations and models' ensem-341 ble mean estimates. We consider spatial scales from 0.1° to 2° , and ensemble size from 342 8 to 128. For all the considered spatial scales, the data-driven models offer precipitation 343 estimates that are significantly more accurate than the CFS reanalysis precipitation prod-344 uct (dashed lines). This highlights the necessity of learning from high-fidelity data (i.e., 345 observations or high-resolution simulations) to represent unresolved processes in climate 346 modeling. Specific to the data-driven models, DiP and CGAN demonstrates similar r347 and RMSE skill, matching or slightly falling behind UNet (solid lines). Meanwhile, CVAE 348 offers optimal r and RMSE skill for spatial scales beyond grid-resolution level (0.1°) . In 349 principle, a supervised learning approach, i.e., UNet, should provide the optimal deter-350 ministic skill. Yet, our results highlight that, for spatial scales that models are not di-351 rectly trained on, a probabilistic model that better exploit the spatial coherency can out-352 perform a supervised learning model. While CVAE has demonstrated this potential, there 353 is room of progress for DiP and CGAN to further improve their deterministic skills. 354

For probabilistic evaluation, we compute the continuous ranked probabilistic skill 355 (CRPS, Fig. 4c), the skill-spread correlation (SSC, Fig. 4d), and the coverage ratio (CR, 356 Fig. 4e) of models' ensemble estimates. For CRPS, the CRPS of a deterministic model, 357 i.e., UNet and CFS reanalysis, is equivalent to the model's mean absolute error. Here, 358 DiP, CGAN, and VAE significantly outperforms UNet and CFS reanalysis. At grid-resolution 359 level, for ensemble size of 8, DiP and CGAN perform similarly, both outperforming CVAE 360 by a large margin. As we gradually double the ensemble size, DiP demonstrates slight 361 advantage over CGAN. This advantage becomes more obvious at larger spatial scales. 362 This result suggests that, compared to CGAN, DiP offers more spatially-coherent prob-363 abilistic estimates. SSC quantifies the reliability of a model's uncertainty estimates: a 364 higher SSC suggests that the model assigns higher/lower forecast uncertainty estimates 365 to forecasts that turn out to have higher/lower biases, which is crucial for decision mak-366 ings. DiP achieves the highest SSC for all spatial scales, followed by CGAN. An increase 367 of ensemble size reduces the statistical error of model's uncertainty estimates, hence in-368 creases model's SSC. This effect is mostly evident for DiP. CR quantifies the ratio that 369 an observation falls into model's ensemble interval, quantifying how well a probabilis-370 tic model is calibrated. Again, DiP achieves the highest CR among the considered mod-371 els, providing a comprehensive range of plausible outcomes. 372

To sum up, DiP verifies competitively compared to alternative data-driven deter-373 ministic/probabilistic approaches, as well as reanalysis precipitation products: for spa-374 tial scales from 0.1° to 2°, DiP matches supervised learning approach in delivering de-375 terministic precipitation estimates (on r and RMSE), and offers optimal probabilistic 376 estimation skills (on CRPS, SSC, and CR). This methodology better meets the READS 377 requirements: it allows us to efficiently generate realistic samples that are faithful to a 378 broad range of resolved circulation schemes, and are diverse to cover most plausible out-379 comes. 380

5 Conclusions

Numerical weather-climate models resolve geophysical fluid dynamics to a finite resolution, necessitating probabilistic inference for unresolved processes. For example, what is the probability that, at millimeter scale, various hydrometeors interact, collide, coalesce to yield precipitation, given circulation status resolved to kilometer scale? If we could accurately and efficiently answer these questions, we could not only better understand, but also better predict the climate.

We follow the data-driven ideology to learn representations of unresolved climate processes from high fidelity data, such as high-resolution simulations and observations.



Figure 4. Performance evaluation using data from 2019 to 2022. The following skill metrics are considered. r: average correlation coefficient between model ensemble mean estimates and observations; RMSE: root mean squared error of model ensemble mean estimate; CRPS: continuous ranked probabilistic score of model ensembles; SSC: spread-skill correlation, where spread is represented using ensemble standard deviation, and skill is represented using model ensemble mean absolute error; CR: coverage ratio, which represents the percentage that grid observation falls into the coverage of ensemble spread. For the probabilistic models, we consider ensemble size from 8 to 128 to compute the skill metrics. All the skill metrics are computed across spatial scales from 0.1° to 2° by spatial pooling.

We point out the limitations of supervised learning approaches in such tasks, and advocate the potential advantages of generative modeling approaches.

To realize these potential advantages, we should steer the learning machine toward verifiable goals of stochastic parameterization, which are quantified in ensemble forecast practices. Hence, based on the requirements of ensemble forecast, we propose the READS (Realism, Efficiency, Adaptability, Diversity, and Sharpness) criteria for probabilistic representation of unresolved climate processes.

To solidify these arguments and provide practical solutions, we consider the problem of numerical precipitation estimation. We develop DiP, a probabilistic diffusion model based methodology to learn stochastic parameterization of precipitation. Compared to existing generative models, DiP approximates a target distribution in a principled, iterative manner, which offers it tremendous fitting capability and controlling flexibility.

Using a Typhoon storm case and four-year evaluation, we demonstrate the advantage of DiP in meeting the READS requirements, as compared to existing data-driven supervised deep learning method (UNet), data-driven probabilistic deep learning method (CVAE and CGAN), as well we moderate/high resolution numerical method (CFS and WRF).

There remain several challenges for our approach to stochastic parameterization. 407 Till now, our model does not provide feedback to the resolved dynamics. It remains to 408 be examined if the learned subgrid-scale noise can trigger circulation regime transitions, 409 and support reliable probabilistic forecast. Also, the ensemble mean estimate from DiP 410 fails to match the performance of CVAE, suggesting room for progress. Finally, to gen-411 erate large ensemble estimates using DiP takes hundreds runs of the deep nets, which 412 brings considerable computation burden in long term simulations. Future works may ex-413 plore diffusion model distillation techniques to accelerate the generation process (Sal-414 imans & Ho, 2022; Song et al., 2023). 415

416 Acknowledgments

⁴¹⁷ This research is supported by National Key R&D Program of China (2021YFA0718000),

⁴¹⁸ National Natural Science Foundation of China (42275174, 42288101) and Chinese Academy

of Science Light of the West Interdisciplinary Research Grant (xbzg-zdsys-202104). We

420 thank Dr. Juanjuan Liu, Dr. Li Dong, Dr. Guiwan Chen, Mr. Jie Chao, and Mr. Yucheng

⁴²¹ Zi for supporting the Human eYe Perceptual Evaluation. The Multi-Source Weighted-

Ensemble Precipitation data are available from https://www.gloh2o.org/mswep/. The

Climate Forecast System Reanalysis data are available from https://climatedataguide.ucar.edu/climate data/climate-forecast-system-reanalysis-cfsr.

425 References

- Beck, H. E., Wood, E. F., Pan, M., Fisher, C. K., Miralles, D. G., Van Dijk, A. I.,
- Adler, R. F. (2019). Mswep v2 global 3-hourly 0.1 precipitation: methodology
 and quantitative assessment. Bulletin of the American Meteorological Society,
 100(3), 473-500.
- ⁴³⁰ Berner, J., Achatz, U., Batte, L., Bengtsson, L., De La Camara, A., Christensen,
- H. M., ... others (2017). Stochastic parameterization: Toward a new view of
 weather and climate models. Bulletin of the American Meteorological Society,
 98(3), 565-588.
- Chen, G., & Wang, W.-C. (2022). Short-term precipitation prediction for con tiguous united states using deep learning. *Geophysical Research Letters*, 49(8),
 e2022GL097904.
- 437 Dorrestijn, J., Crommelin, D. T., Siebesma, A. P., & Jonker, H. J. (2013). Stochas-

438	tic parameterization of shallow cumulus convection estimated from high-resolution
439	Cagne D I Christensen H M Subramanian A C & Monahan A H (2020)
440	Machine learning for stochastic parameterization: Generative adversarial networks
442	in the lorenz'96 model. Journal of Advances in Modeling Earth Systems, 12(3),
443	e2019MS001896.
444	Gneiting, T., & Raftery, A. E. (2007). Strictly proper scoring rules, prediction, and
445	estimation. Journal of the American statistical Association, 102(477), 359–378.
446	Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S.,
447	Bengio, Y. (2014). Generative adversarial nets. Advances in neural informa-
448	tion processing systems, 27.
449	Hardiman, S. C., Dunstone, N. J., Scaife, A. A., Smith, D. M., Comer, R., Nie, Y.,
450	& Ren, HL. (2022). Missing eddy feedback may explain weak signal-to-noise
451	ratios in climate predictions. npj Climate and Atmospheric Science, $5(1), 57$.
452	Harris, L., McRae, A. I., Chantry, M., Dueben, P. D., & Palmer, I. N. (2022). A
453	casts <i>Journal of Advances in Modeling Earth Systems</i> 1/(10) e2022MS003120
454	Ho I Jain A & Abbeel P (2020) Denoising diffusion probabilistic models Ad-
455	vances in neural information processing systems, 33, 6840–6851.
457	Ho, J., & Salimans, T. (2022). Classifier-free diffusion guidance. <i>arXiv preprint</i>
458	arXiv:2207.12598.
459	Holmes, C. C., & Walker, S. G. (2017). Assigning a value to a power likelihood in a
460	general bayesian model. $Biometrika$, $104(2)$, $497-503$.
461	Kingma, D., Salimans, T., Poole, B., & Ho, J. (2021). Variational diffusion models.
462	Advances in neural information processing systems, 34, 21696–21707.
463	Kingma, D. P., & Welling, M. (2013). Auto-encoding variational bayes. <i>arXiv</i>
464	preprint arXiv:1312.6114.
465	Palmer, T. (2019). Stochastic weather and climate models. <i>Nature Reviews Physics</i> ,
466	I(I), 403-4II.
467	Palmer, I. N., Buizza, R., Doblas-Reyes, F., Jung, I., Leutbecher, M., Shutts,
468	G. J., Weishelmer, A. (2009). Stochastic parametrization and model uncer-
470	Pan, B., Anderson, G. J., Goncalves, A., Lucas, D. D., Bonfils, C. J., & Lee, J.
471	(2022). Improving seasonal forecast using probabilistic deep learning. <i>Journal of</i>
472	Advances in Modeling Earth Systems, 14(3), e2021MS002766.
473	Pan, B., Anderson, G. J., Goncalves, A., Lucas, D. D., Bonfils, C. J., Lee, J.,
474	Ma, HY. (2021). Learning to correct climate projection biases. Journal of
475	Advances in Modeling Earth Systems, $13(10)$, e2021MS002509.
476	Pan, B., Hsu, K., AghaKouchak, A., & Sorooshian, S. (2019). Improving precipi-
477	tation estimation using convolutional neural network. Water Resources Research, 55(2), 2201, 2221
478	JJ(J), $2JUI=2J2I$. Dan B. Hay K. Asha-Koushali A. Saraashian S. Is Uisming W. (2010). Determined
479	itation prediction skill for the west coast united states: From short to extended
48U 481	range, Journal of Climate, 32(1), 161–182.
482	Plant, R., & Craig, G. C. (2008). A stochastic parameterization for deep convec-
483	tion based on equilibrium statistics. Journal of the Atmospheric Sciences, 65(1),
484	87–105.
485	Ravuri, S., Lenc, K., Willson, M., Kangin, D., Lam, R., Mirowski, P., others
486	(2021). Skilful precipitation nowcasting using deep generative models of radar.
487	Nature, 597(7878), 672–677.
488	Ronneberger, O., Fischer, P., & Brox, T. (2015). U-net: Convolutional networks
489	tor biomedical image segmentation. In Medical image computing and computer-
490	ussisieu intervention-miccui 2013: 18th international conference, munich, ger- many, october 5-9, 2015, proceedings, part iii 18 (pp. 234–241)
491	many, october 5-3, 2013, proceedings, part ill 10 (pp. 234-241).

- Ruthotto, L., & Haber, E. (2021). An introduction to deep generative modeling.
 GAMM-Mitteilungen, 44 (2), e202100008.
- Saha, S., Nadiga, S., Thiaw, C., Wang, J., Wang, W., Zhang, Q., ... others (2006).
 The ncep climate forecast system. *Journal of Climate*, 19(15), 3483–3517.
- Sakradzija, M., Seifert, A., & Dipankar, A. (2016). A stochastic scale-aware param eterization of shallow cumulus convection across the convective gray zone. Journal
 of Advances in Modeling Earth Systems, 8(2), 786–812.
- Salimans, T., & Ho, J. (2022). Progressive distillation for fast sampling of diffusion
 models. arXiv preprint arXiv:2202.00512.
- Skamarock, W. C., Klemp, J. B., Dudhia, J., Gill, D. O., Liu, Z., Berner, J., ... others (2019). A description of the advanced research wrf version 4. NCAR tech. note ncar/tn-556+ str, 145.
- Sohl-Dickstein, J., Weiss, E., Maheswaranathan, N., & Ganguli, S. (2015). Deep un supervised learning using nonequilibrium thermodynamics. In International con ference on machine learning (pp. 2256–2265).
- ⁵⁰⁷ Song, Y., Dhariwal, P., Chen, M., & Sutskever, I. (2023). Consistency models.
- Song, Y., Sohl-Dickstein, J., Kingma, D. P., Kumar, A., Ermon, S., & Poole, B.
 (2020). Score-based generative modeling through stochastic differential equations.
- arXiv preprint arXiv:2011.13456.
- Stensrud, D. J. (2009). Parameterization schemes: keys to understanding numerical
 weather prediction models. Cambridge University Press.
- Tapiador, F. J., Roca, R., Del Genio, A., Dewitte, B., Petersen, W., & Zhang, F.
 (2019). Is precipitation a good metric for model performance? Bulletin of the American Meteorological Society, 100(2), 223–233.
- Wang, L.-Y., & Tan, Z.-M. (2023). Deep learning parameterization of the tropical cyclone boundary layer. Journal of Advances in Modeling Earth Systems, 15(1), e2022MS003034.
- Yu, S., Hannah, W. M., Peng, L., Bhouri, M. A., Gupta, R., Lin, J., ... oth-
- ers (2023). Climsim: An open large-scale dataset for training high-resolution physics emulators in hybrid multi-scale climate simulators. *arXiv preprint arXiv:2306.08754*.
- ⁵²³ Zhou, S., Gordon, M., Krishna, R., Narcomey, A., Morina, D., & Bernstein, M. S.
- ⁵²⁴ (2019). Hype: human-eye perceptual evaluation of generative models.





0	· –	- 44	<u> </u>			r 0		Г А		~	F 0
0	5 5) [[0 1	5 20	0 23	5 3	iU 3	5 4	-0 4	ら >	×50

	ensemble members	μ	σ	mean absolute error	ur
aDiP (ours)HYPE↑3/5r↑0.571RMSE↓4.96CRPS↓2.45SSC↑0.73CR↑80.5%					
b CGAN HYPE↑ 1/5 r↑ 0.502 RMSE↓ 5.28 CRPS↓ 2.63 SSC↑ 0.58 CR↑ 53.6%					
CVAEHYPE↑0/5r↑0.507RMSE↓4.61CRPS↓2.68SSC↑-0.01CR↑37.9%					

precipitation rate (mm/h)

0	5	10	15	20	25	
D	iP	CGAN	CV	ΆE	Observation	า







Probabilistic diffusion model for stochastic 1 parameterization 2 – a case example of numerical precipitation estimation 3 Baoxiang Pan¹, Leyi Wang², Feng Zhang³, Qingyun Duan⁴, Xin Li⁵, Xiaoduo 4 Pan⁵, Xi Chen¹, Fenghua Ling⁶, Shuguang Wang⁷, Ming Pan⁸, Ziniu Xiao¹ 5 ¹Institute of Atmospheric Physics, Chinese Academy of Sciences 6 ²Chongqing Institute of Big Data, Peking University ³Department of Atmospheric and Oceanic Sciences, Fudan University ⁴National Key Laboratory of Water Disaster Prevention, Hohai University ⁵Institute of Tibetan Plateau Research, Chinese Academy of Sciences ⁶Institute for Climate and Application Research, Nanjing University of Information Science and 8 9 10 11 Technology 12 ⁷School of Atmospheric Sciences, Nanjing University ⁸Scripps Institution of Oceanography, University of California San Diego 13 14

To learn stochastic parameterization, one should steer generative models toward the requirements of ensemble forecasts. We propose criteria of READS (Realism, Efficiency, Adaptability, Diversity, Sharpness) for data-driven stochastic parameterization. In case example of numerical precipitation estimation, probabilistic diffusion model well meets the READS criteria.

Key Points:

15

Corresponding author: Baoxiang Pan, panbaoxiang@lasg.iap.ac.cn

22 Abstract

Estimating the unresolved geophysical processes from resolved geophysical fluid dynam-23 ics is the key for improving numerical weather-climate predictions. While data-driven 24 parameterization for unresolved geophysical processes shows potential, most practices 25 fail to capture the diversity of unresolved geophysical processes that agree with resolved 26 geophysical fluid state. This pitfall undermines the likelihood or severity of simulated 27 weather extremes, and erodes the fidelity of climate projections. We propose the crite-28 ria of READS (Realism, Efficiency, Adaptability, Diversity, Sharpness) for generative mod-29 els to yield reasonable stochastic parameterization. We introduce probabilistic diffusion 30 model, a non-equilibrium thermodynamics inspired deep generative modeling approach, 31 to better meet these criteria. Using a case example of numerical precipitation estima-32 tion, we demonstrate the advantage of the proposed methodology in quickly delivering 33 diverse and faithful estimates for the target unresolved process, as compared to other 34 popular data-driven deterministic and stochastic methods (UNet, variational autoencoder, 35 generative adversarial net), as well as dynamical downscaling method (WRF). We con-36 clude that generative models, in particular, probabilistic diffusion model, can significantly 37 enhance the representation of unresolved geophysical processes in numerical weather-climate 38 predictions. 39

⁴⁰ Plain Language Summary

⁴¹ "Life is a gorgeous robe, crawling with lice", so said Eileen Chang, a Chinese writer
⁴² who enjoyed depicting the awkward discrepancies between ideal and reality. Same metaphor
⁴³ applies to climate models, rooted in physical principles of fluid dynamics and thermo⁴⁴ dynamics, rife with empirics making up the missing components. We use generative AI
⁴⁵ to make up the missing components in climate models, achieving realistic and informa⁴⁶ tive simulations of unresolved climate processes, i.e., precipitation.

47 **1** Introduction

Geophysical fluid dynamics operates across a continuous spectrum of spatiotemporal scales, ranging from micro-scale turbulences to synoptic-scale planetary waves. Their numerical solvers, coming with finite resolution, set a distinction between resolved dynamics and unresolved physical processes, with the latter being approximated as empirical functions of the former. This approximation, known as parameterization, is the source of error in numerical weather and climate predictions (Stensrud, 2009).

Typically, parameterization schemes are deterministic functions, providing a unique 54 tendency accounting for the grid-scale impact of subgrid physical processes in numer-55 ical modeling of geophysical fluid dynamics. However, as we do not explicitly resolve the 56 subgrid physical processes, a probabilistic formulation is advocated (Berner et al., 2017; 57 T. Palmer, 2019): the impact of subgrid physical processes should be described by a prob-58 ability distribution function conditioning on the resolved geophysical fluid dynamics. This 59 probabilistic formulation enables a rigorous and consistent characterization of unresolved 60 physical processes across model resolutions (Sakradzija et al., 2016). Also, it allows subgrid-61 scale noise to trigger crucial circulation regime transitions, supporting reliable proba-62 bilistic forecasts (T. N. Palmer et al., 2009). 63

To make accurate probabilistic representation of unresolved physical processes in numerical weather-climate models, existing efforts proceed along the following three lines. The first, which is straightforward yet lacks theoretical warrant, is to pre-define a perturbation to the parameters, functions, or outputs of deterministic parameterization schemes (Dorrestijn et al., 2013). The second, which is solid yet costly and restricted in scope, is to compute statistics of the equilibrium states of the considered process, via statistical mechanics analysis (Plant & Craig, 2008). The third, which promises remarkable accuracy, efficiency, and flexibility, is to approximate the probability distribution of unresolved physical processes by *learning* from high fidelity data, such as high-resolution simulations or observations (Gagne et al., 2020; Ravuri et al., 2021; Harris et al., 2022).

We proceed along the third line as motivated by the recent breakthrough of prob-74 abilistic machine learning, in particular, probabilistic diffusion models (Sohl-Dickstein 75 et al., 2015; Ho et al., 2020; Song et al., 2020). Probabilistic diffusion models learn to 76 approximate probability distributions in an iterative manner, achieving unprecedented 77 fitting capacity and controlling flexibility in generative modeling tasks. Using a case ex-78 79 ample of numerical precipitation estimation, we specify five requirements for developing data-driven stochastic parameterization schemes. We develop Diffusion based Pre-80 cipitation estimator, dubbed DiP, and demonstrate its unique advantages in meeting these 81 requirements, as compared to existing data-driven deterministic and stochastic param-82 eterization schemes, as well as high resolution dynamical simulation method. 83

⁸⁴ 2 Problem setup and model requirements

We consider a case example of numerical precipitation estimation: given geophys-85 ical fluid dynamics resolved to a finite spatiotemporal resolution, the goal is to estimate 86 the accompanying precipitation process. The challenge lies in that, precipitation results 87 from a complicated chain of processes that are mostly unresolved in numerical models 88 (Tapiador et al., 2019). Any error along this simulation chain may distort the location, 89 timing, or quantity of the precipitation estimate, rendering the estimate useless, even mis-90 leading (Pan, Hsu, AghaKouchak, & Sorooshian, 2019; Pan, Hsu, AghaKouchak, Sorooshian, 91 & Higgins, 2019; Chen & Wang, 2022). 92

Here we consider the region of East and Southeast Asia $(0^{\circ}-40^{\circ}N, 100^{\circ}E-140^{\circ}E)$, 93 where precipitation is driven by diverse circulation regimes. We use the following resolv-94 able dynamical variables to infer precipitation: key primitive variables (meridional and 95 zonal wind velocity, temperature, specific humidity, and geopotential height) at 3 pres-96 sure levels (1000/850/500 hPa), and crucial surface level variables (sea level pressure, 97 surface pressure, surface temperature, and total column precipitable water). These data 98 are obtained by blending observations with short-range weather forecasts to faithfully 99 represent historical circulation states. The data are from the Climate Forecast System 100 Reanalysis project (Saha et al., 2006), coming at spatiotemporal resolution of $0.5^{\circ}/1$ hour 101 for Year 1979-2022. Besides these dynamical variables, we also consider 0.1° elevation 102 data as extra, static predictor. These dynamical and static predictor variables are to-103 gether denoted as \mathbf{x} . Precipitation, as our predict and variable, is denoted by \mathbf{y} . The data 104 are from the Multi-Source Weighted-Ensemble Precipitation product (Beck et al., 2019), 105 which merges gauge, satellite, and reanalysis precipitation records to achieve optimal qual-106 ity. The data come at a $0.1^{\circ}/3$ -hourly resolution for same period. 107

Our objective is to approximate the conditional distribution of $p(\mathbf{y}|\mathbf{x})$, based on favorably large amount of $\{\mathbf{x}, \mathbf{y}\}$ paired data samples. This problem setup differs from a deterministic regression problem setup, which has been widely adopted for learning parameterization schemes (Yu et al., 2023; Wang & Tan, 2023). Specifically, in both deterministic and probabilistic formulations, we design a learning machine Θ , for which the optimal parameter θ^* is obtained by maximizing the overall likelihood of the reference data:

$$\theta^* = \operatorname*{argmax}_{\theta} \sum_{i} \log p_{\theta}(\mathbf{y}_i | \mathbf{x}_i) \tag{1}$$

In a deterministic regression problem setup, given any \mathbf{x} , the learning machine yields the most plausible \mathbf{y} : $\hat{\mathbf{y}} = \underset{\mathbf{y}}{\operatorname{argmax}} \log p_{\theta^*}(\mathbf{y}|\mathbf{x})$. This is often achieved by pre-assuming the distributional form of p_{θ} . For instance, assuming $p_{\theta}(\mathbf{y}|\mathbf{x}) := \mathcal{N}(\mu_{\theta}(\mathbf{x}), \sigma_{\theta}^{2}(\mathbf{x}))$, we have $\hat{\mathbf{y}} = \mu_{\theta^{*}}(\mathbf{x})$, where $\theta^{*} = \underset{\theta}{\operatorname{argmin}} \sum_{i} \left[\frac{(\mu_{\theta}(\mathbf{x}_{i}) - \mathbf{y}_{i})^{2}}{\sigma_{\theta}^{2}(\mathbf{x}_{i})} + \log \sigma_{\theta}^{2}(\mathbf{x}_{i}) \right]$. Such an assumption offers huge computation convenience, yet, it comes with two deficits. First, a predefined distributional form often poorly fits a richly structured target physical process. Second, a deterministic formulation precludes interaction between subgrid noise and resolved dynamics, resulting in biased weather-climate predictions (Hardiman et al., 2022).

In a probabilistic modeling setup, given any \mathbf{x} , the learning machine outputs plau-123 sible y samples: $\hat{\mathbf{y}} \sim p_{\theta^*}(\mathbf{y}|\mathbf{x})$, where p_{θ^*} is a learned distribution subject to no pre-124 defined probability distribution form. This probabilistic formulation allows us to bypass 125 the two deficits of a deterministic formulation, yet, it comes with its own challenges and 126 requirements. To realize the potential, one must steer the learning machine toward ver-127 ifiable goals of stochastic parameterization, which are quantified in ensemble forecast prac-128 tices. We hence suggest the following five criteria for p_{θ^*} based on the requirements of 129 ensemble forecast: 130

131	• Realism : samples from the estimated conditional probability distribution should
132	be indistinguishable from observational samples, regarding either their structure
133	or functionality. This requirement ensures that an accurate probability value can
134	be assigned to the realized observations, either for training or evaluation purposes.
135	Also, the generated samples can fit into the geophysical modeling pipeline, and
136	be as useful as observations for a wide range of subsequent tasks.
137	• Efficiency: a solid approach for developing parameterization schemes is to con-
138	sider each of the possible ways that the subgrid scale process evolve under the grid-
139	scale constraint, to compute the probability of each such "configuration" in the
140	equilibrium ensemble, and generate samples accordingly. This requires excessive
141	human effort and computational resources. Here, we expect a well-trained p_{θ^*} to
142	efficiently generate multiple samples of plausible subgrid physical processes, at least
143	several orders faster than directly resolving the subgrid scale process.
144	• Adaptability: the interaction of subgrid scale physics and large scale dynamics
145	often results in organized weather schemes across scales, ranging from local con-
146	vection to weather fronts. Correspondingly, the model is preferred to automati-
147	cally identify and apply to these organized weather schemes, rather than work-
148	ing at individual computing grids or fixed computing time steps.
149	• Diversity: the estimated conditional probability distribution should cover all plau-
150	sible outcomes, rather than a limited subset of modes. This ensures that all ob-
151	served states are within the <i>cone</i> of model simulations, particularly for extremes.
152	• Sharpness: the estimated conditional probability distribution should generate
153	samples that are faithful to the conditioning information, maximizing the sharp-
154	ness of the simulated distribution. Note that this requirement naturally confronts
155	the diversity requirement: an overly constrained probability estimate may fail to
156	encapsulate observations, which is unreliable; an overly dispersed probability es-
157	timate may lack clear distinction from climatology, which is uninformative. We
158	must carefully balance sharpness and diversity, so that the probability estimate
159	faithfully reflects the intrinsic stochasticity of the considered process.

We coin the term READS by concatenating the initial letters of the five criteria above. Below we introduce DiP, Diffusion based Precipitation estimator, and demonstrate its unique advantage in meeting the READS requirements, as compared to existing deterministic/probabilistic data-driven approaches, as well as high-resolution dynamical simulation approach.

¹⁶⁵ 3 Diffusion based Precipitation estimator (DiP)

3.1 A primer on probabilistic machine learning

The method we develop here falls into the scope of probabilistic machine learning, 167 which applies probability theory to design learning machines that make predictions as 168 probability distributions. Since the target distribution we try to approximate adheres 169 to no pre-defined closed form, a common strategy is to learn a mapping between the tar-170 get distribution and a tractable latent distribution, i.e., standard Gaussian. After learn-171 ing the mapping from optimally large amount of data, we can pass samples from the la-172 tent distribution through the trained model to obtain target distribution samples, hence 173 inferring this target distribution. The key challenge is that, we lack point-to-point cor-174 respondences between samples from the target distribution and the latent distribution, 175 hence lacking straightforward supervision signals to enable learning (Ruthotto & Haber, 176 2021). A popular solution is to build bijective mapping between the target distribution 177 and the latent distribution, therefore establishing correspondences. Probabilistic diffu-178 sion models excel in this task by establishing the bijection in an iterative manner. Be-179 low we outline how this is achieved. Mathematical and implementation details are given 180 in Supporting Information S1. 181

¹⁸² 3.2 Basics

166

¹⁸³ Diffusion model approximates a target distribution $p(\mathbf{y})$ by reversing a Gaussian ¹⁸⁴ process (Fig. 1): the forward Gaussian process turns $p(\mathbf{y})$ into standard Gaussian $\mathcal{N}(\mathbf{0}, \mathbf{I})$ ¹⁸⁵ (Fig. 1a, Eq. 2); we learn to iteratively reverse this Gaussian process, mapping $\mathcal{N}(\mathbf{0}, \mathbf{I})$ ¹⁸⁶ to $p(\mathbf{y})$ (Fig. 1b-d), hence achieving generative modeling. Following D. Kingma et al. (2021), ¹⁸⁷ the forward Gaussian process is pre-defined as:

$$p(\mathbf{z}_t | \mathbf{y}) := \mathcal{N}(\alpha_t \mathbf{y}, \sigma_t^2 \mathbf{I})$$
(2)

Here \mathbf{z}_t is latent variable indexed by $t \in [0, 1]$, α_t / σ_t is monotonically decreasing/increasing function of t, strictly bounded by [0, 1]. Eq. 2 therefore bridges $p(\mathbf{y}) = p(\mathbf{z}_0|\mathbf{y})$ and $\mathcal{N}(\mathbf{0}, \mathbf{I}) = p(\mathbf{z}_1|\mathbf{y})$ (Fig. 1a). We reverse Eq. 2 to turn $\mathcal{N}(\mathbf{0}, \mathbf{I})$ into $p(\mathbf{y})$, using a chain of variational distributions (Fig. 1b):

$$p_{\theta}(\mathbf{z}_{t_{i-1}}|\mathbf{z}_{t_i}) = \mathcal{N}(\mu_{\theta}(\mathbf{z}_{t_i}), \mathbf{\Sigma}_{\theta}(\mathbf{z}_{t_i})), \quad i \in [1, T]$$
(3)

Here $0 = t_0 < t_1 < t_2 < ... < t_T = 1$ is arbitrary discretization of time; $\{\mu_{\theta}, \Sigma_{\theta}\}$ are

learnable mean vector and covariance matrix, trained by maximizing the overall data like lihood (Supporting Information S1.1):

$$\log p_{\theta}(\mathbf{y}) = \mathbb{E}_{p(\mathbf{z}_{0}|\mathbf{y})} \log p(\mathbf{y}|\mathbf{z}_{0}) - D_{\mathrm{KL}} \left(p(\mathbf{z}_{1}|\mathbf{y}) || p(\mathbf{z}_{1}) \right) - \sum_{i=1}^{T} \mathbb{E}_{p(\mathbf{z}_{t_{i}}|\mathbf{y})} D_{\mathrm{KL}} \left(p(\mathbf{z}_{t_{i-1}}|\mathbf{z}_{t_{i}},\mathbf{y}) || p_{\theta}(\mathbf{z}_{t_{i-1}}|\mathbf{z}_{t_{i}}) \right)$$
(4)

Given Eq. 2, to maximize Eq. 4 is approximately equivalent to minimizing the Fisher divergence between the data and model distributions (Supporting Information S1.2):

$$\theta^* = \operatorname*{argmax}_{\theta} \log p_{\theta}(\mathbf{y}) \approx \operatorname*{argmin}_{\theta} \sum_{i=1}^{T} \mathbb{E}_{p(\mathbf{z}_{t_i}|\mathbf{y})} \left\| \nabla \log p(\mathbf{z}_{t_i}|\mathbf{y}) - \epsilon_{\mathrm{NN}_{\theta}}(\mathbf{z}_{t_i}) \right\|_2$$
(5)

Here $\epsilon_{NN_{\theta}}$ is a neural network parameterization of $\nabla \log p(\mathbf{z}_{t_i}|\mathbf{y})$, known as the score func-

tion. Based on the learned score estimates, we can derive $p_{\theta}(\mathbf{z}_{t_{i-1}}|\mathbf{z}_{t_i}) = \mathcal{N}(\mu_{\theta}(\mathbf{z}_{t_i}), \boldsymbol{\Sigma}_{\theta}(\mathbf{z}_{t_i}))$ (Supporting Information S1.2) and sample it, starting with $p(\mathbf{z}_1) = \mathcal{N}(\mathbf{0}, \mathbf{I})$, ending with

 $n(\mathbf{z}_0) \approx n(\mathbf{y})$

198 $p(\mathbf{z}_0) \approx p(\mathbf{y}).$



Figure 1. Overview of diffusion model. We map target distribution (synoptic-scale precipitation field, **a** left) to a same dimensional standard Gaussian distribution (**a** right) through a pre-defined Gaussian process (**a** bottom, Eq. 2). Color denotes probability distribution function value for an individual precipitation field pixel (here we select the center pixel) through diffusion time t = [0, 1], lines show the diffusion trajectories of individual pixels for randomly selected samples, matrix plots show the noisified precipitation field (sample of Typhoon Lekima, 0000 UTC 09 August 2019, centered at 26.5°N, 114.4°E) across diffusion time (**a** top). We approximate the target distribution by reversing the Gaussian process, using a series of variational distributions (**b**, Eq. 3), which are trained by maximizing the data likelihood (Eq. 4-5). We include conditioning information to approximate conditional distribution of a same Typhoon event (**c**). We apply *classifier-free guidance* to control the impact of the conditioning information versus the latent variable in explaining the variability of the target variable for the same event (**d**, Eq. 6). By enhancing the guidance strength ω , we suppress the variance of the resulting conditional probability distribution (**c**/**d** right). The plots are supported by logarithm transformed precipitation observational data for Year 2019, and the trained diffusion models.

¹⁹⁹ 3.3 Conditioning

To generate \mathbf{y} samples that are faithful to the conditioning information \mathbf{x} , we need 200 to approximate the conditional distribution $p(\mathbf{y}|\mathbf{x})$. To achieve this, we include \mathbf{x} dur-201 ing training and sampling (Fig. 1c-d). A direct inclusion of \mathbf{x} does not specify the im-202 pact of **x** versus \mathbf{z}_t in explaining the variability of **y** (Fig. 1c, Holmes & Walker 2017). 203 To tackle the potential misspecification, and having \mathbf{x} effectively control the learned dis-204 tribution, we resort to *classifier-free guidance* (Ho & Salimans, 2022, Fig. 1d): we learn 205 two sets of neural networks: $\epsilon_{\rm NN}(\mathbf{z}_{t_i}) / \epsilon_{\rm NN}^c(\mathbf{z}_{t_i}, \mathbf{x})$, so to approximate the unconditional/conditional 206 scores: $\nabla \log p(\mathbf{z}_{t_i}|\mathbf{y}) / \nabla \log p(\mathbf{z}_{t_i}|\mathbf{y}, \mathbf{x})$. Based upon these two sets of score estimates, 207 we compose score estimators for synthetic distributions $p_{\omega}(\mathbf{z}_{t_i}|\mathbf{x},\mathbf{y}) \propto p(\mathbf{x}|\mathbf{y},\mathbf{z}_{t_i})^{\omega} p(\mathbf{z}_{t_i}|\mathbf{y})$: 208

$$\nabla \log p_{\omega}(\mathbf{z}_{t_i}|\mathbf{y}, \mathbf{x}) = \nabla \log p(\mathbf{z}_{t_i}|\mathbf{y}) + \omega \nabla \log p(\mathbf{x}|\mathbf{y}, \mathbf{z}_{t_i})$$

= $\nabla \log p(\mathbf{z}_{t_i}|\mathbf{y}) + \omega (\nabla \log p(\mathbf{z}_{t_i}|\mathbf{x}, \mathbf{y}) - \nabla \log p(\mathbf{z}_{t_i}|\mathbf{y}))$ (6)
 $\approx \epsilon_{\mathrm{NN}}(\mathbf{z}_{t_i}) + \omega (\epsilon_{\mathrm{NN}}^c(\mathbf{z}_{t_i}, \mathbf{x}) - \epsilon_{\mathrm{NN}}(\mathbf{z}_{t_i}))$

Here ω is guidance scale coefficient, balancing the diversity and sharpness of the learned conditional distribution:

 $\omega = 1$: assuming impact of **x** has been perfectly accounted by $\epsilon_{\text{NN}}^c(\mathbf{z}_{t_i}, \mathbf{x})$ (Fig. 1c).

 $\omega < 1$: suppressing impact of **x**, pervading the distribution toward climatology.

 $\omega > 1$: raising impact of **x**, sharpening the distribution toward more likely values (Fig. 1d).

We now apply score estimates of $p_{\omega}(\mathbf{z}_{t_i}|\mathbf{x}, \mathbf{y})$ to sample $p(\mathbf{y}|\mathbf{x})$, following a same strategy described in Sec. 3.2. The value of ω is empirically determined based on the probabilistic forecasting skill of its resulting model.

3.4 Baselines and implementation details

We compare the DiP methodology with popular deterministic and stochastic data-

- driven methods and moderate/high resolution dynamical simulation method, including:
- **UNet**: a de-facto choice for image-to-image regression tasks, using neural network consisting symmetric convolution and deconvolution blocks (Ronneberger et al., 2015).
- Conditional variational autoencoder (CVAE): a probabilistic deep learn ing method that maximizes a lower bound of data likelihood to learn latent variable model for a target conditional distribution (D. P. Kingma & Welling, 2013;
 Pan et al., 2022).
- Conditional generative adversarial net (CGAN): a probabilistic deep learning method in which a generative network learns to approximate a target conditional distribution, under the guidance of a discriminative network that distinguishes generated samples and true samples (Goodfellow et al., 2014; Pan et al., 2021; Ravuri et al., 2021).
- CFS reanalysis precipitation product (CFSR): an optimized combination of CMAP (CPC Merged Analysis of Precipitation), daily gauge observations, and CFS background 6-hourly precipitation analysis (Saha et al., 2006).
- Dynamical downscaling using WRF: refining coarsely resolved climate pro cesses via high resolution numerical geophysical fluid dynamics solver and accom panying parameterization schemes, using Advanced Research Version 4.2 of Weather
 Research and Forecasting (WRF-ARW V4.2, Skamarock et al. 2019).
- For all the data-driven models, including DiP, we use data from 1979-2016/2017-2018/2019-2022 for training/validation/test. Considering the computation cost and the

characteristic scale of atmospheric dynamics, all the data-driven models operate at a synoptic scale $(8^{\circ} \times 8^{\circ})$: we randomly crop paired predictor and predictand field data within the study region for model training. The model structures, hyper-parameter setups, and training details are given in Supporting Information S2.

²⁴⁵ **3.5 Evaluation**

We verify models' performances using a suite of skill metrics corresponding to the 246 READS criteria. We apply Human eYe Perceptual Evaluation (HYPE, Zhou et al. 2019) 247 and power spectrum analysis to determine models' sample fidelity. We use Pearson cor-248 relation coefficient (r) and Root Mean Squared Error (RMSE) between observations and 249 models' ensemble mean estimations to quantify models' deterministic prediction skills. 250 We apply Continuous Ranked Probabilistic Skill (CRPS) to measure the accuracy of the 251 predicted probabilities and the sharpness of the forecast distribution. We compute model's 252 skill spread correlation (SSC) to quantify the reliability of a model's uncertainty esti-253 mates. We compute the ratio that observations falls into model's ensemble intervals (CR). 254 We record the computing time of the considered models. All the skill metrics are com-255 puted across spatial scales from 0.1° to 2° by aggregating neighbourhood grids. For de-256 tails, see Supporting Information S3. 257

258 4 Results

259 4.1 Case study

We start with a case example to compare models' performances. We consider the 260 storm process associated with Typhoon Lekima, which ranks as the third costliest ty-261 phoon in Chinese history. We show $8^{\circ} \times 8^{\circ}$ observed and simulated precipitation rate 262 maps along the typhoon trajectory (Fig. 2). Here, observations (Fig. 2a) present a clear 263 ring structure of intense precipitation surrounding the typhoon eye before landing (0000 264 UTC 04 August 2019 - 0000 UTC 08 August 2019), with maximum precipitation rate 265 reaching 100 mm/h. The eyewall structure gradually dissipates through two landings (1800 266 UTC 09 August and 1200 UTC 11 August), leaving a tightly curved rainband wrapping 267 into a relatively well-defined centre. 268

The large-scale patterns of precipitation estimates from the data-driven models (Fig. 2b-269 e) and CFS reanalysis (Fig. 2f) roughly agree with observations (Fig. 2a), due to a shared 270 circulation constraint from CFS reanalysis. For WRF dynamical downscaling (Fig. 2g), 271 despite careful spectral nudging, the results do not strictly follow the observed typhoon 272 trajectory, particularly after landing (1800 UTC 09 August). This is due to the chaotic 273 nature of geophysical fluid dynamics. The fine-scale structure differs significantly among 274 models: DiP (Fig. 2b) produces the most realistic small-scale details, creating a clear eve-275 wall structure and associated spiral rainband, with intense precipitation matching ob-276 servations at relatively correct locations. CGAN (Fig. 2c) can generate intense precip-277 itations surrounding the typhoon eye. Yet, the estimates come with poor spatial struc-278 ture, with neighboring grids loosely correlated, and the rainband barely depictable. CVAE 279 (Fig. 2d) and UNet (Fig. 2e) offer similar, blurry estimates, failing to distinct charac-280 teristic typhoon eyewall and rainband structures. Besides, both models miss precipita-281 tion extremes, with maximum precipitation estimates below 30 mm/h. CFS reanalysis 282 (Fig. 2f) shares similar drawbacks as CVAE and UNet, largely due to biases from the 283 assimilated data sources and errors from precipitation related model parameterization 284 schemes. WRF simulation (Fig. 2g) makes overly confined, extremely intense (approx-285 imately 150 mm/h) precipitation estimates, following the finely resolved, yet potentially 286 misaligned circulation state estimates. 287

We further inspect the probabilistic models (DiP, CGAN, and CVAE) through the lens of the READS requirements (Sec. 2). For an individual snapshot of precipitation



Figure 2. Observed and simulated $8^{\circ} \times 8^{\circ}$ precipitation rate maps along the trajectory of Typhoon Lekima, from 0000 UTC 04 August 2019 to 0000 UTC 12 August 2019. a: precipitation observations from MSWEP. b-d: randomly selected samples of ensemble precipitation estimates using DiP/CGAN/CVAE. e: deterministic precipitation estimates using UNet. f: CFS reanalysis precipitation with resolution of 0.2° . g: precipitation estimates using WRF dynamical simulation, with resolution of ~ 3 km. The typhoon trajectory from WRF simulation considerably diverges from observations after the first landing (1800 UTC 09 August). For after landing results, we show precipitation rate maps surrounding WRF simulated typhoon center.

estimate centering around 22.7°N, 125.9°E at 0000 UTC 06 August 2019, we show models' ensemble members, ensemble mean and standard deviation, ensemble mean absolute error, as well as radial/orientation averaged power spectrum (Fig. 3). We compute
a suite of skill metrics corresponding to the READS requirements.

• **Realism**: we measure human climate experts' error rate in detecting observation 294 from model estimates: for DiP/CGAN/CVAE, 3/1/0 out of 5 climate scientist eval-295 uators fail to detect the observation from 15 randomly generated model estimates, 296 suggesting the optimal spatial coherency of DiP estimates. Additionally, we in-297 spect the spatial structure of precipitation estimates by computing their average 298 spectrum power as function of spatial frequency and orientation: DiP and CGAN 299 well reproduce the spatial variability across spatial scales and orientations. Mean-300 while, WRF significantly overestimates spatial variability; CVAE, UNet and CFSR 301 significantly underestimate spatial variability for high spatial frequency and all 302 orientations. 303

- Efficiency: all the probabilistic models demonstrate advantageous efficiency compared to high-resolution numerical simulation: DiP/CGAN/CVAE generate 100 member ensemble estimates of 0.1° precipitation field within approximately 100/2/2
 seconds on a NVIDIA GeForce RTX 4090 GPU. Here, DiP is two-orders slower
 than CGAN and CVAE due to its iterative generation nature. As a comparison,
 a deterministic WRF simulation takes around 5 hours in a 32-core CPU machine.
- Adaptability: data-driven models are often reported to struggle with extremes, due to unreasonable learning objective setups, as well as approximation, optimization, and statistical errors. While the typhoon case we consider here is featured by extreme precipitation, DiP successfully reproduces the maximum precipitation rate and characteristic typhoon rainfall structures, suggesting its adaptability for extreme cases. We further report models' performances for various weather schemes in Sec. 4.2.
- **Diversity-Sharpness tradeoff**: we measure the diversity of models' ensemble 317 estimates by computing the percentage that a grid point observation falls into model's 318 ensemble interval. Here, 80.5%/53.6%/29.7% grid point observations are within 319 the 16-member ensemble interval from DiP/CGAN/CVAE. Grid points where ob-320 servations fall above/below the ensemble interval are stippled with red/black. These 321 results suggest the peculiar advantage of DiP in delivering broad range of plau-322 sible outcomes. We further investigate model's sharpness subject to a "proper" 323 level of diversity. By "proper", we mean that the probability estimate accurately 324 reflects the intrinsic stochasticity of the considered process, which is not directly 325 measurable and requires statistical inference. A good indicator is how model's en-326 semble spread aligns with model's skill. DiP achieves the highest spread-skill cor-327 relation, assigning high/low forecast uncertainty estimates to predictions with high/low 328 errors. We further consider the spatial correlation between the ensemble mean es-329 timate and observation, as well as the mean absolute error between each ensem-330 ble member and observation. The high skill values of DiP suggest that its ensem-331 ble dispersion centers around observation, requiring no ensemble pruning. Finally, 332 we report models' continuous ranked probability scores, which considers both pre-333 diction diversity and sharpness. DiP achieves the optimal performance under this 334 proper scoring rule (Gneiting & Raftery, 2007). 335
 - 4.2 Skill evaluation

336

We evaluated models' overall performances using test set data from 2019 to 2022. We report a suite of deterministic and probabilistic skill metrics for the considered models in Fig. 4.



Figure 3. Precipitation estimates centering around 22.7° N, 125.9° E at 0000 UTC 06 August 2019, using DiP (a), CGAN (b), and CVAE (c). The columns show models' ensemble members, ensemble mean, ensemble standard deviation, ensemble mean absolute error, grid points where observation is not encapsulated by ensemble spread (red/black stipple for under/over estimation, background colored based on observation), and radial/orientation averaged power spectrum for observation and all the considered models, including DiP, CGAN, CVAE, UNet, CFS reanalysis, and WRF. The following skill metrics are computed. HYPE: human climate experts' error rate in detecting observation from model estimates; r: spatial correlation between model ensemble mean estimate; CRPS: continuous ranked probabilistic score of model ensembles; SSC: spread-skill correlation, where spread is represented using ensemble standard deviation, and skill is represented using model ensemble mean absolute error; CR: coverage ratio, which represent the percentage that grid observation falls into the coverage of ensemble spread.

For deterministic evaluation, we compute the correlation coefficient (r, Fig. 4a) and 340 the root mean squared error (RMSE, Fig. 4b) between observations and models' ensem-341 ble mean estimates. We consider spatial scales from 0.1° to 2° , and ensemble size from 342 8 to 128. For all the considered spatial scales, the data-driven models offer precipitation 343 estimates that are significantly more accurate than the CFS reanalysis precipitation prod-344 uct (dashed lines). This highlights the necessity of learning from high-fidelity data (i.e., 345 observations or high-resolution simulations) to represent unresolved processes in climate 346 modeling. Specific to the data-driven models, DiP and CGAN demonstrates similar r347 and RMSE skill, matching or slightly falling behind UNet (solid lines). Meanwhile, CVAE 348 offers optimal r and RMSE skill for spatial scales beyond grid-resolution level (0.1°) . In 349 principle, a supervised learning approach, i.e., UNet, should provide the optimal deter-350 ministic skill. Yet, our results highlight that, for spatial scales that models are not di-351 rectly trained on, a probabilistic model that better exploit the spatial coherency can out-352 perform a supervised learning model. While CVAE has demonstrated this potential, there 353 is room of progress for DiP and CGAN to further improve their deterministic skills. 354

For probabilistic evaluation, we compute the continuous ranked probabilistic skill 355 (CRPS, Fig. 4c), the skill-spread correlation (SSC, Fig. 4d), and the coverage ratio (CR, 356 Fig. 4e) of models' ensemble estimates. For CRPS, the CRPS of a deterministic model, 357 i.e., UNet and CFS reanalysis, is equivalent to the model's mean absolute error. Here, 358 DiP, CGAN, and VAE significantly outperforms UNet and CFS reanalysis. At grid-resolution 359 level, for ensemble size of 8, DiP and CGAN perform similarly, both outperforming CVAE 360 by a large margin. As we gradually double the ensemble size, DiP demonstrates slight 361 advantage over CGAN. This advantage becomes more obvious at larger spatial scales. 362 This result suggests that, compared to CGAN, DiP offers more spatially-coherent prob-363 abilistic estimates. SSC quantifies the reliability of a model's uncertainty estimates: a 364 higher SSC suggests that the model assigns higher/lower forecast uncertainty estimates 365 to forecasts that turn out to have higher/lower biases, which is crucial for decision mak-366 ings. DiP achieves the highest SSC for all spatial scales, followed by CGAN. An increase 367 of ensemble size reduces the statistical error of model's uncertainty estimates, hence in-368 creases model's SSC. This effect is mostly evident for DiP. CR quantifies the ratio that 369 an observation falls into model's ensemble interval, quantifying how well a probabilis-370 tic model is calibrated. Again, DiP achieves the highest CR among the considered mod-371 els, providing a comprehensive range of plausible outcomes. 372

To sum up, DiP verifies competitively compared to alternative data-driven deter-373 ministic/probabilistic approaches, as well as reanalysis precipitation products: for spa-374 tial scales from 0.1° to 2°, DiP matches supervised learning approach in delivering de-375 terministic precipitation estimates (on r and RMSE), and offers optimal probabilistic 376 estimation skills (on CRPS, SSC, and CR). This methodology better meets the READS 377 requirements: it allows us to efficiently generate realistic samples that are faithful to a 378 broad range of resolved circulation schemes, and are diverse to cover most plausible out-379 comes. 380

5 Conclusions

Numerical weather-climate models resolve geophysical fluid dynamics to a finite resolution, necessitating probabilistic inference for unresolved processes. For example, what is the probability that, at millimeter scale, various hydrometeors interact, collide, coalesce to yield precipitation, given circulation status resolved to kilometer scale? If we could accurately and efficiently answer these questions, we could not only better understand, but also better predict the climate.

We follow the data-driven ideology to learn representations of unresolved climate processes from high fidelity data, such as high-resolution simulations and observations.



Figure 4. Performance evaluation using data from 2019 to 2022. The following skill metrics are considered. r: average correlation coefficient between model ensemble mean estimates and observations; RMSE: root mean squared error of model ensemble mean estimate; CRPS: continuous ranked probabilistic score of model ensembles; SSC: spread-skill correlation, where spread is represented using ensemble standard deviation, and skill is represented using model ensemble mean absolute error; CR: coverage ratio, which represents the percentage that grid observation falls into the coverage of ensemble spread. For the probabilistic models, we consider ensemble size from 8 to 128 to compute the skill metrics. All the skill metrics are computed across spatial scales from 0.1° to 2° by spatial pooling.

We point out the limitations of supervised learning approaches in such tasks, and advocate the potential advantages of generative modeling approaches.

To realize these potential advantages, we should steer the learning machine toward verifiable goals of stochastic parameterization, which are quantified in ensemble forecast practices. Hence, based on the requirements of ensemble forecast, we propose the READS (Realism, Efficiency, Adaptability, Diversity, and Sharpness) criteria for probabilistic representation of unresolved climate processes.

To solidify these arguments and provide practical solutions, we consider the problem of numerical precipitation estimation. We develop DiP, a probabilistic diffusion model based methodology to learn stochastic parameterization of precipitation. Compared to existing generative models, DiP approximates a target distribution in a principled, iterative manner, which offers it tremendous fitting capability and controlling flexibility.

Using a Typhoon storm case and four-year evaluation, we demonstrate the advantage of DiP in meeting the READS requirements, as compared to existing data-driven supervised deep learning method (UNet), data-driven probabilistic deep learning method (CVAE and CGAN), as well we moderate/high resolution numerical method (CFS and WRF).

There remain several challenges for our approach to stochastic parameterization. 407 Till now, our model does not provide feedback to the resolved dynamics. It remains to 408 be examined if the learned subgrid-scale noise can trigger circulation regime transitions, 409 and support reliable probabilistic forecast. Also, the ensemble mean estimate from DiP 410 fails to match the performance of CVAE, suggesting room for progress. Finally, to gen-411 erate large ensemble estimates using DiP takes hundreds runs of the deep nets, which 412 brings considerable computation burden in long term simulations. Future works may ex-413 plore diffusion model distillation techniques to accelerate the generation process (Sal-414 imans & Ho, 2022; Song et al., 2023). 415

416 Acknowledgments

⁴¹⁷ This research is supported by National Key R&D Program of China (2021YFA0718000),

⁴¹⁸ National Natural Science Foundation of China (42275174, 42288101) and Chinese Academy

of Science Light of the West Interdisciplinary Research Grant (xbzg-zdsys-202104). We

420 thank Dr. Juanjuan Liu, Dr. Li Dong, Dr. Guiwan Chen, Mr. Jie Chao, and Mr. Yucheng

⁴²¹ Zi for supporting the Human eYe Perceptual Evaluation. The Multi-Source Weighted-

Ensemble Precipitation data are available from https://www.gloh2o.org/mswep/. The

Climate Forecast System Reanalysis data are available from https://climatedataguide.ucar.edu/climate data/climate-forecast-system-reanalysis-cfsr.

425 References

- Beck, H. E., Wood, E. F., Pan, M., Fisher, C. K., Miralles, D. G., Van Dijk, A. I.,
- Adler, R. F. (2019). Mswep v2 global 3-hourly 0.1 precipitation: methodology
 and quantitative assessment. Bulletin of the American Meteorological Society,
 100(3), 473-500.
- ⁴³⁰ Berner, J., Achatz, U., Batte, L., Bengtsson, L., De La Camara, A., Christensen,
- H. M., ... others (2017). Stochastic parameterization: Toward a new view of
 weather and climate models. Bulletin of the American Meteorological Society,
 98(3), 565-588.
- Chen, G., & Wang, W.-C. (2022). Short-term precipitation prediction for con tiguous united states using deep learning. *Geophysical Research Letters*, 49(8),
 e2022GL097904.
- 437 Dorrestijn, J., Crommelin, D. T., Siebesma, A. P., & Jonker, H. J. (2013). Stochas-

438	tic parameterization of shallow cumulus convection estimated from high-resolution
439	Cagne D I Christensen H M Subramanian A C & Monahan A H (2020)
440	Machine learning for stochastic parameterization: Generative adversarial networks
442	in the lorenz'96 model. Journal of Advances in Modeling Earth Systems, 12(3),
443	e2019MS001896.
444	Gneiting, T., & Raftery, A. E. (2007). Strictly proper scoring rules, prediction, and
445	estimation. Journal of the American statistical Association, 102(477), 359–378.
446	Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S.,
447	Bengio, Y. (2014). Generative adversarial nets. Advances in neural informa-
448	tion processing systems, 27.
449	Hardiman, S. C., Dunstone, N. J., Scaife, A. A., Smith, D. M., Comer, R., Nie, Y.,
450	& Ren, HL. (2022). Missing eddy feedback may explain weak signal-to-noise
451	ratios in climate predictions. npj Climate and Atmospheric Science, $5(1), 57$.
452	Harris, L., McRae, A. I., Chantry, M., Dueben, P. D., & Palmer, I. N. (2022). A
453	casts <i>Journal of Advances in Modeling Earth Systems</i> 1/(10) e2022MS003120
454	Ho I Jain A & Abbeel P (2020) Denoising diffusion probabilistic models Ad-
455	vances in neural information processing systems, 33, 6840–6851.
457	Ho, J., & Salimans, T. (2022). Classifier-free diffusion guidance. <i>arXiv preprint</i>
458	arXiv:2207.12598.
459	Holmes, C. C., & Walker, S. G. (2017). Assigning a value to a power likelihood in a
460	general bayesian model. $Biometrika$, $104(2)$, $497-503$.
461	Kingma, D., Salimans, T., Poole, B., & Ho, J. (2021). Variational diffusion models.
462	Advances in neural information processing systems, 34, 21696–21707.
463	Kingma, D. P., & Welling, M. (2013). Auto-encoding variational bayes. <i>arXiv</i>
464	preprint arXiv:1312.6114.
465	Palmer, T. (2019). Stochastic weather and climate models. <i>Nature Reviews Physics</i> ,
466	I(I), 403-4II.
467	Palmer, I. N., Buizza, R., Doblas-Reyes, F., Jung, I., Leutbecher, M., Shutts,
468	G. J., Weishelmer, A. (2009). Stochastic parametrization and model uncer-
470	Pan, B., Anderson, G. J., Goncalves, A., Lucas, D. D., Bonfils, C. J., & Lee, J.
471	(2022). Improving seasonal forecast using probabilistic deep learning. <i>Journal of</i>
472	Advances in Modeling Earth Systems, 14(3), e2021MS002766.
473	Pan, B., Anderson, G. J., Goncalves, A., Lucas, D. D., Bonfils, C. J., Lee, J.,
474	Ma, HY. (2021). Learning to correct climate projection biases. Journal of
475	Advances in Modeling Earth Systems, $13(10)$, e2021MS002509.
476	Pan, B., Hsu, K., AghaKouchak, A., & Sorooshian, S. (2019). Improving precipi-
477	tation estimation using convolutional neural network. Water Resources Research, 55(2), 2201, 2221
478	JJ(J), 2JUI=2J2I. Dan B. Hey K. Asha-Koushalt, A. Sorooshian S. & Uissing W. (2010). Determined
479	itation prediction skill for the west coast united states: From short to extended
48U 481	range, Journal of Climate, 32(1), 161–182.
482	Plant, R., & Craig, G. C. (2008). A stochastic parameterization for deep convec-
483	tion based on equilibrium statistics. Journal of the Atmospheric Sciences, 65(1),
484	87–105.
485	Ravuri, S., Lenc, K., Willson, M., Kangin, D., Lam, R., Mirowski, P., others
486	(2021). Skilful precipitation nowcasting using deep generative models of radar.
487	Nature, 597(7878), 672–677.
488	Ronneberger, O., Fischer, P., & Brox, T. (2015). U-net: Convolutional networks
489	tor biomedical image segmentation. In Medical image computing and computer-
490	ussisieu intervention-miccui 2013: 18th international conference, munich, ger- many, october 5-9, 2015, proceedings, part iii 18 (pp. 234–241)
491	many, october 5-3, 2013, proceedings, part ill 10 (pp. 234-241).

- Ruthotto, L., & Haber, E. (2021). An introduction to deep generative modeling.
 GAMM-Mitteilungen, 44 (2), e202100008.
- Saha, S., Nadiga, S., Thiaw, C., Wang, J., Wang, W., Zhang, Q., ... others (2006).
 The ncep climate forecast system. *Journal of Climate*, 19(15), 3483–3517.
- Sakradzija, M., Seifert, A., & Dipankar, A. (2016). A stochastic scale-aware param eterization of shallow cumulus convection across the convective gray zone. Journal
 of Advances in Modeling Earth Systems, 8(2), 786–812.
- Salimans, T., & Ho, J. (2022). Progressive distillation for fast sampling of diffusion
 models. arXiv preprint arXiv:2202.00512.
- Skamarock, W. C., Klemp, J. B., Dudhia, J., Gill, D. O., Liu, Z., Berner, J., ... others (2019). A description of the advanced research wrf version 4. NCAR tech. note ncar/tn-556+ str, 145.
- Sohl-Dickstein, J., Weiss, E., Maheswaranathan, N., & Ganguli, S. (2015). Deep un supervised learning using nonequilibrium thermodynamics. In International con ference on machine learning (pp. 2256–2265).
- ⁵⁰⁷ Song, Y., Dhariwal, P., Chen, M., & Sutskever, I. (2023). Consistency models.
- Song, Y., Sohl-Dickstein, J., Kingma, D. P., Kumar, A., Ermon, S., & Poole, B.
 (2020). Score-based generative modeling through stochastic differential equations.
- arXiv preprint arXiv:2011.13456.
- Stensrud, D. J. (2009). Parameterization schemes: keys to understanding numerical
 weather prediction models. Cambridge University Press.
- Tapiador, F. J., Roca, R., Del Genio, A., Dewitte, B., Petersen, W., & Zhang, F.
 (2019). Is precipitation a good metric for model performance? Bulletin of the American Meteorological Society, 100(2), 223–233.
- Wang, L.-Y., & Tan, Z.-M. (2023). Deep learning parameterization of the tropical cyclone boundary layer. Journal of Advances in Modeling Earth Systems, 15(1), e2022MS003034.
- Yu, S., Hannah, W. M., Peng, L., Bhouri, M. A., Gupta, R., Lin, J., ... oth-
- ers (2023). Climsim: An open large-scale dataset for training high-resolution physics emulators in hybrid multi-scale climate simulators. *arXiv preprint arXiv:2306.08754*.
- ⁵²³ Zhou, S., Gordon, M., Krishna, R., Narcomey, A., Morina, D., & Bernstein, M. S.
- ⁵²⁴ (2019). Hype: human-eye perceptual evaluation of generative models.

Supporting Information for "Probabilistic diffusion model for stochastic parameterization – a case example of numerical precipitation estimation"

Baoxiang Pan¹, Leyi Wang², Feng Zhang³, Qingyun Duan⁴, Xin Li⁵, Xiaoduo

Pan⁵, Xi Chen¹, Fenghua Ling⁶, Shuguang Wang⁷, Ming Pan⁸, Ziniu Xiao¹

¹Institute of Atmospheric Physics, Chinese Academy of Sciences

²Chongqing Institute of Big Data, Peking University

³Department of Atmospheric and Oceanic Sciences, Fudan University

⁴National Key Laboratory of Water Disaster Prevention, Hohai University

 $^5 \mathrm{Institute}$ of Tibetan Plateau Research, Chinese Academy of Sciences

⁶Institute for Climate and Application Research, Nanjing University of Information Science and Technology

⁷School of Atmospheric Sciences, Nanjing University

⁸Scripps Institution of Oceanography, University of California San Diego

Contents of this file

- S1. Details of probabilistic diffusion model
- S2. Baseline models
- S3. Evaluation metrics

S1. Details of probabilistic diffusion model

The theory and practice of probabilistic diffusion models can be math-heavy and convoluted. We provide mathematical and implementation details of the deployed probabilistic diffusion model. For a friendly tutorial, see Luo (2022). For more information and useful learning materials, see Sohl-Dickstein et al. (2015), Ho et al. (2020), Kingma et al. (2021), and Song et al. (2020).

:

S1.1 Decomposition of $\log p(y)$ using latent $z_{0:1}$

Diffusion model is an explicit likelihood based generative model. Its learning objective function is a factorization of data likelihood defined over latent variables $\mathbf{z}_{0:1} = \{\mathbf{z}_{t_0=0}, \mathbf{z}_{t_1}, \mathbf{z}_{t_2}, ..., \mathbf{z}_{t_T=1}\}$. This factorization stands at the core of variational view of diffusion models. A step-by-step derivation is given below.

First, for a pre-defined $p(\mathbf{z}_{0:1}|\mathbf{y})$, we have:

$$\log p(\mathbf{y}) = E_{p(\mathbf{z}_{0:1}|\mathbf{y})} \left[\log p(\mathbf{y})\right] = E_{p(\mathbf{z}_{0:1}|\mathbf{y})} \left[\log \frac{p(\mathbf{z}_{0:1}, \mathbf{y})}{p(\mathbf{z}_{0:1}|\mathbf{y})}\right]$$
(1)

Given $p(\mathbf{z}_t | \mathbf{y}) := \mathcal{N}(\alpha_t \mathbf{y}, \sigma_t^2 \mathbf{I}), 0 = t_0 < t_1 < t_2 < \dots < t_T = 1$, we have:

$$p(\mathbf{z}_{0:1}, \mathbf{y}) = p(\mathbf{y} | \mathbf{z}_{0:1}) p(\mathbf{z}_{t_0} | \mathbf{z}_{t_1:t_T}) p(\mathbf{z}_{t_1} | \mathbf{z}_{t_2:t_T}) \dots p(\mathbf{z}_{t_{T-1}} | \mathbf{z}_{t_T}) p(\mathbf{z}_{t_T})$$

$$= p(\mathbf{y} | \mathbf{z}_{t_0}) p(\mathbf{z}_{t_0} | \mathbf{z}_{t_1}) p(\mathbf{z}_{t_1} | \mathbf{z}_{t_2}) \dots p(\mathbf{z}_{t_{T-1}} | \mathbf{z}_{t_T}) p(\mathbf{z}_{t_T})$$

$$= p(\mathbf{y} | \mathbf{z}_{t_0}) p(\mathbf{z}_{t_T}) \prod_{i=1}^{T} p(\mathbf{z}_{t_{i-1}} | \mathbf{z}_{t_i})$$
(2)

and

$$p(\mathbf{z}_{0:1}|\mathbf{y}) = p(\mathbf{z}_{t_0}|\mathbf{y})p(\mathbf{z}_{t_1}|\mathbf{z}_{t_0},\mathbf{y})p(\mathbf{z}_{t_2}|\mathbf{z}_{t_{0:t_1}},\mathbf{y})...p(\mathbf{z}_{t_T}|\mathbf{z}_{t_{0:t_{T-1}}},\mathbf{y})$$

$$= p(\mathbf{z}_{t_0}|\mathbf{y})p(\mathbf{z}_{t_1}|\mathbf{z}_{t_0})p(\mathbf{z}_{t_2}|\mathbf{z}_{t_1})...p(\mathbf{z}_{t_T}|\mathbf{z}_{t_{T-1}})$$

$$= p(\mathbf{z}_{t_0}|\mathbf{y})\prod_{i=1}^{T} p(\mathbf{z}_{t_i}|\mathbf{z}_{t_{i-1}})$$
(3)

Plug Eq. 2 and Eq. 3 into Eq. 1, we have:

$$\log p(\mathbf{y}) = E_{p(\mathbf{z}_{0:1}|\mathbf{y})} \left[\log \frac{p(\mathbf{z}_{0:1}, \mathbf{y})}{p(\mathbf{z}_{0:1}|\mathbf{y})} \right]$$
$$= E_{p(\mathbf{z}_{0:1}|\mathbf{y})} \left[\log \frac{p(\mathbf{y}|\mathbf{z}_{t_0})p(\mathbf{z}_{t_T}) \prod_{i=1}^{T} p(\mathbf{z}_{t_{i-1}}|\mathbf{z}_{t_i})}{p(\mathbf{z}_{t_0}|\mathbf{y}) \prod_{i=1}^{T} p(\mathbf{z}_{t_i}|\mathbf{z}_{t_{i-1}})} \right]$$
$$= E_{p(\mathbf{z}_{0:1}|\mathbf{y})} \left[\log \frac{p(\mathbf{y}|\mathbf{z}_{t_0})p(\mathbf{z}_{t_T})}{p(\mathbf{z}_{t_0}|\mathbf{y})} + \sum_{i=1}^{T} \log \frac{p(\mathbf{z}_{t_{i-1}}|\mathbf{z}_{t_i})}{p(\mathbf{z}_{t_i}|\mathbf{z}_{t_{i-1}})} \right]$$
(4)

While Eq. 4 can be decomposed into differentiable terms, it involves estimating expectation over two random variables: $\{\mathbf{z}_{t_{i-1}}, \mathbf{z}_{t_i}\}$, which may have high variance (Luo, 2022). To achieve a robust estimate, we re-write $p(\mathbf{z}_{t_i}|\mathbf{z}_{t_{i-1}})$ as:

:

$$p(\mathbf{z}_{t_i}|\mathbf{z}_{t_{i-1}}) = p(\mathbf{z}_{t_i}|\mathbf{z}_{t_{i-1}}, \mathbf{y}) = \frac{p(\mathbf{z}_{t_{i-1}}|\mathbf{z}_{t_i}, \mathbf{y})p(\mathbf{z}_{t_i}|\mathbf{y})}{p(\mathbf{z}_{t_{i-1}}|\mathbf{y})}$$
(5)

Plug Eq. 5 into Eq. 4, we have:

$$\begin{split} \log p(\mathbf{y}) &= E_{p(\mathbf{z}_{0:1}|\mathbf{y})} \left[\log \frac{p(\mathbf{y}|\mathbf{z}_{t_{0}})p(\mathbf{z}_{t_{T}})}{p(\mathbf{z}_{t_{0}}|\mathbf{y})} + \sum_{i=1}^{T} \log \frac{p(\mathbf{z}_{t_{i-1}}|\mathbf{z}_{t_{i}})}{p(\mathbf{z}_{t_{i-1}}|\mathbf{z}_{i})} \right] \\ &= E_{p(\mathbf{z}_{0:1}|\mathbf{y})} \left[\log \frac{p(\mathbf{y}|\mathbf{z}_{t_{0}})p(\mathbf{z}_{t_{T}})}{p(\mathbf{z}_{t_{0}}|\mathbf{y})} + \sum_{i=1}^{T} \log \frac{p(\mathbf{z}_{t_{i-1}}|\mathbf{z}_{t_{i}})}{p(\mathbf{z}_{t_{i-1}}|\mathbf{x}_{i})} \right] \\ &= E_{p(\mathbf{z}_{0:1}|\mathbf{y})} \left[\log \frac{p(\mathbf{y}|\mathbf{z}_{t_{0}})p(\mathbf{z}_{t_{T}})}{p(\mathbf{z}_{t_{0}}|\mathbf{y})} + \sum_{i=1}^{T} \log \frac{p(\mathbf{z}_{t_{i-1}}|\mathbf{z}_{t_{i}})}{p(\mathbf{z}_{t_{i-1}}|\mathbf{x}_{i})} + \sum_{i=1}^{T} \log \frac{p(\mathbf{z}_{t_{i-1}}|\mathbf{y})}{p(\mathbf{z}_{t_{i-1}}|\mathbf{y})} \right] \\ &= E_{p(\mathbf{z}_{0:1}|\mathbf{y})} \left[\log \frac{p(\mathbf{y}|\mathbf{z}_{0})p(\mathbf{z}_{t_{T}})}{p(\mathbf{z}_{t_{0}}|\mathbf{y})} + \sum_{i=1}^{T} \log \frac{p(\mathbf{z}_{t_{i-1}}|\mathbf{z}_{t_{i}})}{p(\mathbf{z}_{t_{i-1}}|\mathbf{z}_{t_{i}})} + \log \frac{p(\mathbf{z}_{t_{0}}|\mathbf{y})}{p(\mathbf{z}_{t_{T}}|\mathbf{y})} \right] \\ &= E_{p(\mathbf{z}_{0:1}|\mathbf{y})} \left[\log p(\mathbf{y}|\mathbf{z}_{0}) + \sum_{i=1}^{T} \log \frac{p(\mathbf{z}_{t_{i-1}}|\mathbf{z}_{t_{i}})}{p(\mathbf{z}_{t_{i-1}}|\mathbf{z}_{t_{i}},\mathbf{y})} + \log \frac{p(\mathbf{z}_{t_{T}})}{p(\mathbf{z}_{t_{T}}|\mathbf{y})} \right] \\ &= E_{p(\mathbf{z}_{0:1}|\mathbf{y})} \left[\log p(\mathbf{y}|\mathbf{z}_{0}) + E_{p(\mathbf{z}_{T}|\mathbf{y})} \log \frac{p(\mathbf{z}_{t_{T}})}{p(\mathbf{z}_{t_{1}}|\mathbf{z}_{t_{i}},\mathbf{y})} + \log \frac{p(\mathbf{z}_{t_{T}})}{p(\mathbf{z}_{t_{T}}|\mathbf{y})} \right] \\ &= E_{p(\mathbf{z}_{0:1}|\mathbf{y})} \log p(\mathbf{y}|\mathbf{z}_{0}) + E_{p(\mathbf{z}_{T}|\mathbf{y})} \log \frac{p(\mathbf{z}_{t_{T}})}{p(\mathbf{z}_{t_{T}}|\mathbf{y})} + \sum_{i=1}^{T} E_{p(\mathbf{z}_{i}|\mathbf{y})} \log \frac{p(\mathbf{z}_{t_{i-1}}|\mathbf{z}_{t_{i}},\mathbf{y})}{p(\mathbf{z}_{t_{T}}|\mathbf{y})} \right] \\ &= E_{p(\mathbf{z}_{0}|\mathbf{y})} \log p(\mathbf{y}|\mathbf{z}_{0}) + E_{p(\mathbf{z}_{T}|\mathbf{y})} \log \frac{p(\mathbf{z}_{t_{T}})}{p(\mathbf{z}_{t_{T}}|\mathbf{y})} + \sum_{i=1}^{T} E_{p(\mathbf{z}_{i}|\mathbf{y})} \log \frac{p(\mathbf{z}_{t_{i-1}}|\mathbf{z}_{t_{i}},\mathbf{y})}{p(\mathbf{z}_{t_{i-1}}|\mathbf{z}_{t_{i}},\mathbf{y})} \\ &= \mathbb{E}_{p(\mathbf{z}_{0}|\mathbf{y})} \log p(\mathbf{y}|\mathbf{z}_{0}) - D_{\mathrm{KL}}(p(\mathbf{z}_{t_{T}}|\mathbf{y})||p(\mathbf{z}_{t_{T}})) - \sum_{i=1}^{T} \mathbb{E}_{p(\mathbf{z}_{i}|\mathbf{y})} D_{\mathrm{KL}}(p(\mathbf{z}_{t_{i-1}}|\mathbf{z}_{t_{i}},\mathbf{y})||p(\mathbf{z}_{t_{i-1}}|\mathbf{z}_{t_{i}})) \\ &= \mathbb{E}_{p(\mathbf{z}_{0}|\mathbf{y})} \log p(\mathbf{y}|\mathbf{z}_{0}) - D_{\mathrm{KL}}(p(\mathbf{z}_{1}|\mathbf{y})||p(\mathbf{z}_{1})) - \sum_{i=1}^{T} \mathbb{E}_{p(\mathbf{z}_{i}|\mathbf{y})} D_{\mathrm{KL}}(p(\mathbf{z}_{t_{i-1}}|\mathbf{z}_{i},\mathbf{y})||p(\mathbf{z}_{t_{i-1}}|\mathbf{z}_{i})) \end{aligned}$$

We now take a close examination of the three terms on the right side of Eq. 6:

•
$$\mathbb{E}_{p(\mathbf{z}_0|\mathbf{y})} \log p(\mathbf{y}|\mathbf{z}_0)$$
:] given that $\alpha_0 \coloneqq 1, \sigma_0 \coloneqq 0$, we have $p(\mathbf{z}_0|\mathbf{y}) = \begin{cases} \delta, \mathbf{z}_0 = y \\ 0, \text{else} \end{cases}$, and $\begin{cases} \delta, y = \mathbf{z}_0 \end{cases}$

$$p(\mathbf{y}|\mathbf{z}_0) = \begin{cases} \delta, y = \mathbf{z}_0 \\ 0, \text{else} \end{cases}, \ \delta \text{ is Dirac function. Thus, we have } \mathbb{E}_{p(\mathbf{z}_0|\mathbf{y})} \log p(\mathbf{y}|\mathbf{z}_0) = 0. \\ \text{November 23, 2023, 5:12am} \end{cases}$$

Х - З

• $D_{\mathrm{KL}}(p(\mathbf{z}_1|\mathbf{y})||p(\mathbf{z}_1))$: given that $\alpha_1 \coloneqq 0, \sigma_1 \coloneqq 1$, we have $p(\mathbf{z}_1|\mathbf{y}) \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, which is agnostic of \mathbf{y} , this leads to $p(\mathbf{z}_1) \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, therefore, $D_{\mathrm{KL}}(p(\mathbf{z}_1|\mathbf{y})||p(\mathbf{z}_1)) = 0$.

:

• $\sum_{i=1}^{T} \mathbb{E}_{p(\mathbf{z}_{t_i}|\mathbf{y})} D_{\mathrm{KL}}(p(\mathbf{z}_{t_{i-1}}|\mathbf{z}_{t_i},\mathbf{y})||p(\mathbf{z}_{t_{i-1}}|\mathbf{z}_{t_i})$: for any *i*, we can derive an analytical Gaussian form of $p(\mathbf{z}_{t_{i-1}}|\mathbf{z}_{t_i},\mathbf{y})$ (see below). For fine enough discretization of time, $p(\mathbf{z}_{t_{i-1}}|\mathbf{z}_{t_i})$ is also Gaussian, which is represented as a variational distribution parameterized by neural networks. Thus, we obtain an analytical form of $D_{\mathrm{KL}}(p(\mathbf{z}_{t_{i-1}}|\mathbf{z}_{t_i},\mathbf{y})||p(\mathbf{z}_{t_{i-1}}|\mathbf{z}_{t_i})$, suitable for stochastic gradient optimization.

Given the analysis above, to maximize $\log p(\mathbf{y})$ is approximately equivalent to minimizing the following time averaging Kullback–Leibler divergence term:

$$\log p(\mathbf{y}) \approx -\sum_{i=1}^{T} \mathbb{E}_{p(\mathbf{z}_{t_i}|\mathbf{y})} D_{\mathrm{KL}} \left(p(\mathbf{z}_{t_{i-1}}|\mathbf{z}_{t_i}, \mathbf{y}) || p(\mathbf{z}_{t_{i-1}}|\mathbf{z}_{t_i}) \right)$$
(7)

Below we derive analytical form of $p(\mathbf{z}_{t_{i-1}}|\mathbf{z}_{t_i},\mathbf{y})$ and relate it with score estimate of $p(\mathbf{z}_{t_i}|\mathbf{y})$, which enables robust optimization and high-quality generative modeling.

S1.2 Analytical/parameterized form of $p(\mathbf{z}_{t_{i-1}}|\mathbf{z}_{t_i},\mathbf{y})/p(\mathbf{z}_{t_i}|\mathbf{z}_{t_{i-1}})$

To derive the analytical form of $p(\mathbf{z}_{t_{i-1}}|\mathbf{z}_{t_i},\mathbf{y})$, we have:

$$p(\mathbf{z}_{t_{i-1}}|\mathbf{z}_{t_{i}},\mathbf{y}) = p(\mathbf{z}_{t_{i}}|\mathbf{z}_{t_{i-1}},\mathbf{y}) \frac{p(\mathbf{z}_{t_{i-1}}|\mathbf{y})}{p(\mathbf{z}_{t_{i}}|\mathbf{y})} = p(\mathbf{z}_{t_{i}}|\mathbf{z}_{t_{i-1}}) \frac{p(\mathbf{z}_{t_{i-1}}|\mathbf{y})}{p(\mathbf{z}_{t_{i}}|\mathbf{y})}$$

$$\propto \exp\left(-\frac{1}{2}\left(\frac{(\mathbf{z}_{t_{i}} - \frac{\alpha_{t_{i}}}{\alpha_{t_{i-1}}}\mathbf{z}_{t_{i-1}})^{2}}{\sigma_{t_{i}}^{2} - \frac{\alpha_{t_{i}}^{2}}{\alpha_{t_{i-1}}^{2}}\sigma_{t_{i-1}}^{2}} + \frac{(\mathbf{z}_{t_{i-1}} - \alpha_{t_{i-1}}\mathbf{y})^{2}}{\sigma_{t_{i-1}}^{2}} - \frac{(\mathbf{z}_{t_{i}} - \alpha_{t_{i}}\mathbf{y})^{2}}{\sigma_{t_{i}}^{2}}\right)\right)$$

$$= \exp\left(-\frac{1}{2}\left(\frac{1}{\sigma_{t_{i-1}}^{2}(1 - \frac{\alpha_{t_{i}}^{2}\sigma_{t_{i-1}}^{2}}{\alpha_{t_{i-1}}^{2}\sigma_{t_{i}}^{2}})}\mathbf{z}_{t_{i-1}}^{2} - 2\left(\frac{\frac{\alpha_{t_{i}}}{\alpha_{t_{i-1}}}\mathbf{z}_{t_{i}}}{\sigma_{t_{i-1}}^{2}\sigma_{t_{i-1}}^{2}}} + \frac{\alpha_{t_{i-1}}\mathbf{y}}{\sigma_{t_{i-1}}^{2}}\right)\mathbf{z}_{t_{i-1}} + C(\mathbf{z}_{t_{i}},\mathbf{y})\right)\right)$$
(8)

Hence $p(\mathbf{z}_{t_{i-1}}|\mathbf{z}_{t_i},\mathbf{y}) = \mathcal{N}(\tilde{\mu}_{t_i}, \tilde{\boldsymbol{\Sigma}}_{t_i})$, where:

$$\tilde{\mu}_{t_{i}} = \frac{\alpha_{t_{i-1}}}{\alpha_{t_{i}}} \mathbf{z}_{t_{i}} - \frac{\alpha_{t_{i-1}}}{\alpha_{t_{i}}} (\sigma_{t_{i}}^{2} - \frac{\alpha_{t_{i}}^{2}}{\alpha_{t_{i-1}}^{2}} \sigma_{t_{i-1}}^{2}) \frac{\mathbf{z}_{t_{i}} - \alpha_{t_{i}} \mathbf{y}}{\sigma_{t_{i}}^{2}}$$

$$= \frac{\alpha_{t_{i-1}}}{\alpha_{t_{i}}} \mathbf{z}_{t_{i}} + \frac{\alpha_{t_{i-1}}}{\alpha_{t_{i}}} (\sigma_{t_{i}}^{2} - \frac{\alpha_{t_{i}}^{2}}{\alpha_{t_{i-1}}^{2}} \sigma_{t_{i-1}}^{2}) \nabla \log p(\mathbf{z}_{t_{i}} | \mathbf{y})$$
(9)

and

$$\tilde{\Sigma}_{t_i} = \sigma_{t_{i-1}}^2 \left(1 - \frac{\sigma_{t_{i-1}}^2}{\sigma_{t_i}^2} \frac{\alpha_{t_i}^2}{\alpha_{t_{i-1}}^2}\right) \mathbf{I}$$
(10)

November 23, 2023, 5:12am

Given fine enough discretization of time, $p(\mathbf{z}_{t_{i-1}}|\mathbf{z}_{t_i})$ is also Gaussian, i.e. $p(\mathbf{z}_{t_{i-1}}|\mathbf{z}_{t_i}) = \mathcal{N}(\mu_{t_i}, \mathbf{\Sigma}_{t_i})$. We parameterize μ_{t_i} in accordance with the analytical form of $\tilde{\mu}_{t_i}$:

:

$$\mu_{t_i} = \frac{\alpha_{t_{i-1}}}{\alpha_{t_i}} \mathbf{z}_{t_i} + \frac{\alpha_{t_{i-1}}}{\alpha_{t_i}} (\sigma_{t_i}^2 - \frac{\alpha_{t_i}^2}{\alpha_{t_{i-1}}^2} \sigma_{t_{i-1}}^2) \epsilon_{\text{NN}}(\mathbf{z}_{t_i})$$
(11)

where $\epsilon_{\text{NN}}(\mathbf{z}_{t_i})$ is neural network parameterization of $\nabla \log p(\mathbf{z}_{t_i}|\mathbf{y})$. Following Ho & Salimans (2022), we consider the following simplied representation of Σ_{t_i} :

$$\Sigma_{t_i} = \frac{\sigma_{t_i}^{2v_{t_i}}}{\sigma_{t_{i-1}}^{2v_{t_i}}} \tilde{\Sigma}_{t_i}$$
(12)

where $v_{t_i} \coloneqq 0.5$ for all time steps.

Given the analytical/parameterized form of $p(\mathbf{z}_{t_{i-1}}|\mathbf{z}_{t_i},\mathbf{y})/p(\mathbf{z}_{t_i}|\mathbf{z}_{t_{i-1}})$, we have:

$$\begin{aligned} \epsilon_{\mathrm{NN}}^{*}(\mathbf{z}_{t_{i}}) &= \underset{\epsilon_{\mathrm{NN}}(\mathbf{z}_{t_{i}})}{\operatorname{argmin}} \mathbb{E}_{p(\mathbf{z}_{t_{i}}|\mathbf{y})} D_{\mathrm{KL}} \left(p(\mathbf{z}_{t_{i-1}}|\mathbf{z}_{t_{i}},\mathbf{y}) || p(\mathbf{z}_{t_{i-1}}|\mathbf{z}_{t_{i}}) \right) \\ &\approx \underset{\epsilon_{\mathrm{NN}}(\mathbf{z}_{t_{i}})}{\operatorname{argmin}} \mathbb{E}_{p(\mathbf{z}_{t_{i}}|\mathbf{y})} \left[(\tilde{\mu}_{t_{i}} - \mu_{t_{i}})^{T} \Sigma_{t_{i}}^{-1} (\tilde{\mu}_{t_{i}} - \mu_{t_{i}}) \right] \\ &= \underset{\epsilon_{\mathrm{NN}}(\mathbf{z}_{t_{i}})}{\operatorname{argmin}} \mathbb{E}_{p(\mathbf{z}_{t_{i}}|\mathbf{y})} \left[\tilde{\Sigma}_{t_{i}} \frac{\sigma_{t_{i}}^{2(1-v_{t_{i}})}}{\sigma_{t_{i-1}}^{2(1-v_{t_{i}})}} \left\| \nabla \log p(\mathbf{z}_{t_{i}}|\mathbf{y}) - \epsilon_{\mathrm{NN}}(\mathbf{z}_{t_{i}}) \right\|_{2} \right] \\ &\approx \underset{\epsilon_{\mathrm{NN}}(\mathbf{z}_{t_{i}})}{\operatorname{argmin}} \mathbb{E}_{p(\mathbf{z}_{t_{i}}|\mathbf{y})} \left\| \nabla \log p(\mathbf{z}_{t_{i}}|\mathbf{y}) - \epsilon_{\mathrm{NN}}(\mathbf{z}_{t_{i}}) \right\|_{2} \end{aligned}$$

We apply standard stochastic gradient descent to obtain $\epsilon_{NN}^*(\mathbf{z}_{t_i})$. Thereafter, we can approximate $\{\mu_{t_i}, \boldsymbol{\Sigma}_{t_i}\}$ using Eq. 11 and 12. Based on $p_{\theta}(\mathbf{z}_{t_{i-1}}|\mathbf{z}_{t_i})$, we carry out iterative ancestral sampling: note that $p(\mathbf{z}_1) = \mathcal{N}(\mathbf{0}, \mathbf{I})$ given the forward Gaussian Process setup. Therefore, starting from standard Gaussian samples, we iteratively generate samples of $\mathbf{z}_{t_{T-1}}, \mathbf{z}_{t_{T-2}}, ..., \mathbf{z}_{t_0}$, using the learned distributions of $p(\mathbf{z}_{t_{i-1}}|\mathbf{z}_{t_i}), i = T, T-1, ..., 1$. Finally, $p(\mathbf{y}|\mathbf{z}_0) = \begin{cases} \delta, \quad y = \mathbf{z}_0 \\ 0, \quad \text{else} \end{cases}$.

S1.3 $\{\alpha_t, \sigma_t\}$: noise schedule

In diffusion model, we define a Gaussian process to map target distribution to a standard Gaussian. $\{\alpha_t, \sigma_t\}$ specifies the noise schedule, quantifying how fast the target distribution is diminished through the diffusion process. A better noise schedule allows efficient

November 23, 2023, 5:12am

optimization and improved model likelihood, which may be achieved via optimization (Kingma et al., 2021). Here, for simplicity, we adopt a pre-defined noise schedule, following Ho & Salimans (2022) and Nichol & Dhariwal (2021). Specifically, we parameterize $\{\alpha_t, \sigma_t\}$ as a function of λ_t :

:

$$\alpha_t^2 = 1 - \sigma_t^2 = \frac{1}{1 + e^{-\lambda_t}}$$
(14)

where:

$$\lambda_t = -2\log\tan(at+b) \tag{15}$$

Here $b = \arctan(e^{-\frac{\lambda_{\max}}{2}})$, $a = \arctan(e^{-\frac{\lambda_{\min}}{2}}) - b$, t is uniformly sampled from [0, 1]. { $\lambda_{\min} = -20, \lambda_{\max} = 20$ } are hyper-parameters. This noise schedule represents a hyperbolic secant distribution modified to be supported on a bounded interval (Ho & Salimans, 2022).

S1.4 Encoding of diffusion time step

Diffusion model is an iterative generative model, involving a hierarchy of neural network models $\epsilon_{\text{NN}}(\mathbf{z}_t)$ to approximate score functions $\nabla \log p(\mathbf{z}_t | \mathbf{y})$ at multiple noise levels. While this hierarchy of neural network models can be learned separately, in practice, we often adopt a time-dependent neural network, using an vector embedding of t to account for the impact of learning objective difference for different noise levels. Following Song et al. (2020), we incorporate the time information via Gaussian random features, i.e.: embedding $(t) = [\sin(2\pi\omega t); \cos(2\pi\omega t)]$, where $\omega \sim \mathcal{N}(\mathbf{0}, s\mathbf{I})$, s = 1 is a pre-defined scaling parameter.

S1.5 Model architecture and training details

The neural networks we apply for unconditional/conditional score estimates are timedependent UNets with structures illustrated in Fig. S1 and Fig. S2. For now we do not include attention mechanism for computation efficiency. Both models take into input of noisified precipitation field and nosification scale, and outputs the score estimate. We

X - 6

embed the time information, and stack the time embedding as extra channel to all UNet blocks. Each contract block consists of a long chain of $\{C_{3\times3} + N + \text{ReLU}\}_3$, and a short chain of $\{C_{1\times1}\}_1$, concatenated as a residual block, $C_{n\times n}$ is convolution layer with kernel receptive field of size $n \times n$, N is group normalization, ReLU is rectified linear unit function. Each expand block consists of a long chain of $\{R_2+C_{3\times3}+N+\text{ReLU}\}_3$, and a short chain of $\{R_2, C_{1\times1}\}_1$, concatenating as a residual block, R_n resize the data by n times using linear interpolation. We start with channel size of 128, and double/shrink the channel size by 2 along each contract/expand block. For the conditional score estimating neural network, we includes the conditioning information. This conditioning information is deterministic precipitation estimation, offered by a separate UNnet that takes into input of dynamical field information. In this sense, the conditional score estimating neural network tries to recover and add details of the precipitation information discarded by the deterministic precipitation estimator.

:

We use data from 1979-2016/2017-2018/2019-2022 for training/validation/test. We keep same data splitting strategy for all data-driven models considered in this study. To train the unconditional model, we randomly crop precipitation field data of size 80×80 ($8^{\circ} \times 8^{\circ}$), add random scale noise to the data, and use the unconditional diffusion model to estimate the score. We use ADAM optimizer and an initial learning rate of 10^{-3} . We halve the learning rate if validation loss is not decreasing for 10 epochs. To train the conditional model, we include conditioning information from a UNet based deterministic precipitation estimate, the rest settings are same as the unconditional case.

S2. Baseline models

S2.1 UNet

We consider UNet as 1) a deterministic baseline and 2) the conditioning information extractor for all the generative models. UNet a unique convolutional neural network architecture suited for image-relevant tasks. Here, the model takes into input of resolved dynamical field information and static elevation information, and outputs a deterministic precipitation field estimate. The dynamical field information is provided by a 9-hour (including 3 previous/current/future hours), $8^{\circ} \times 8^{\circ}$ circulation field data, with 19 channels representing 19 dynamical variables, including key primitive variables (meridional and zonal wind velocity, temperature, specific humidity, and geopotential height) at 3 pressure levels (1000/850/500 hPa), and crucial surface level variables (sea level pressure, surface pressure, surface temperature, and total column precipitable water). This dynamical field information is first pre-processed through 3D convolution blocks, and concatenated with preprocessed elevation information, before feeding into a 2D UNet. The UNet applies a convolution based contracting path to capture precipitation relevant dynamical field information, and a symmetric transposed convolution based expanding path to gradually refine precipitation field estimates. Skip connections between symmetrical convolution and transposed convolution blocks are applied to force deeper neural network layers to learn meaningful representations that are not well captured by shallower layers. The learning objective is to minimize the squared error between estimated and observed precipitation. Underlying this objective function is the assumption that $p(\mathbf{y}|\mathbf{x})$ is Gaussian, with identical error covariance for any \mathbf{x} and any grid point. See Fig. S3 for UNet model architecture.

:

S2.2 Conditional variational autoencoder (CVAE)

Conditional variational autoencoder (CVAE) is deep neural network powered probabilistic graphical model. To learn a non-linear latent variable model for the target conditional distribution $p(\mathbf{y}|\mathbf{x})$, CVAE constructs a bijective mapping between $p(\mathbf{y}|\mathbf{x})$ and a tractable latent distribution $p(\mathbf{z}|\mathbf{x})$, using an encoder-decoder neural network architecture. The encoder q_{ϕ} approximates $p(\mathbf{z}|\mathbf{x}, \mathbf{y})$ as a variational Gaussian distribution; the decoder p_{ψ} approximates $p(\mathbf{y}|\mathbf{x}, \mathbf{z})$ using the conditioning information \mathbf{x} and the learned latent vector \mathbf{z} . To approximate the target conditional distribution, $\{q_{\phi}, p_{\psi}\}$ are jointly trained to maximize the following evidence lower bound (ELBO) of the data log likelihood:

:

$$\text{ELBO} = \mathbb{E}_{\mathbf{z} \sim q_{\phi}} \log p_{\psi} - \beta D_{\text{KL}} (q_{\phi} \| p(\mathbf{z} | \mathbf{x}))$$
(16)

Here $p(\mathbf{z}|\mathbf{x})$ is assumed to be standard Gaussian; β is a parameter balancing sample diversity and sample accordance to the conditioning information, similar to the functionality of ω in diffusion model. To train CVAE, we run mini-batches of $\{\mathbf{x}, \mathbf{y}\}$ samples through $\{q_{\phi}, p_{\psi}\}$, and update their parameters to maximize the ELBO, using stochastic gradient ascent. To generate novel samples of $p(\mathbf{y}|\mathbf{x})$, we draw \mathbf{z} samples from $p(\mathbf{z}|\mathbf{x})$ and pass them together with \mathbf{x} through the optimal p_{ψ}^* . See Fig. S4 for model architecture details.

S2.3 Conditional generative adversarial net (CGAN)

Conditional generative adversarial net (CGAN) approximates a target conditional distribution $p(\mathbf{y}|\mathbf{x})$ by setting up a "game" between two neural networks. The generator network G takes into input of the conditioning information \mathbf{x} and random noise \mathbf{z} to create samples that are intended to come from the target distribution; the discriminator network D is a binary classifier, optimized to differentiate between generated samples and true samples:

$$L_D = \mathbb{E}_{\mathbf{y}} (D(\mathbf{y})) - \mathbb{E}_{\mathbf{z}} (D(G(\mathbf{x}, \mathbf{z})))$$
(17)

November 23, 2023, 5:12am

The generator network is optimized to 1) fool the discriminator network, and 2) draw the mean of generated samples close to ground truth observations ():

:

$$L_G = \mathbb{E}_{\mathbf{z}} \Big(D \big(G(\mathbf{x}, \mathbf{z}) \big) \Big) - \lambda \mathbb{E}_{\{\mathbf{y}, \mathbf{z}\}} || G(\mathbf{x}, \mathbf{z}) - \mathbf{y} ||_2$$
(18)

Here λ is a parameter that balances sample diversity and sample accordance to the conditioning information, similar to the functionality of ω/β in diffusion/CVAE model. To train CGAN, we run mini-batches of $\{\mathbf{x}, \mathbf{y}, \mathbf{z}\}$ samples through $\{G, D\}$, and apply stochastic gradient ascent to maximize L_D and L_G . We keep the optimal G^* if it offers best skill performance (measured by continuous ranked probabilistic skill score, as it is a proper scoring rule, see below) for the validation set. To generate novel samples of $p(\mathbf{y}|\mathbf{x})$, we draw \mathbf{z} samples from random noise and pass them together with \mathbf{x} through the optimal G^* . See Fig. S5 for model architecture details.

S2.4 Dynamical downscaling using WRF

We include comparison to a dynamical simulation approach for numerical precipitation estimation. Here the Advanced Research Version 4.2 of Weather Research and Forecasting (WRF-ARW V4.2, Skamarock et al. 2019) is deployed for simulation of a Typhoon precipitation case (Typhoon Lekima, 0000 UTC 04 August 2019 -0000 UTC 12 August 2019). WRF-ARW refines the coarsely resolved climate processes at regional scale, using high-resolution numerical geophysical fluid dynamics solver and a suite of accompanying parameterization schemes. We apply Global Forecast System reanalysis data to provide the initial and boundary condition for the considered precipitation cases. We apply spectral nudging of wind for the outer domain to ensure consistency between the simulated large-scale circulations and the analysis fields. The simulated domains are delineated in Fig. S6. The selected parameterization schemes as listed in Tab. S1.

S3. Evaluation metrics

S3.1 Human eYe Perceptual Evaluation (HYPE)

We apply a simplified Human eYe Perceptual Evaluation (HYPE) to assess the sample quality of models' precipitation estimates, relying on human climate scientists' and climate model end-users' perceptions. We measure human climate experts' error rate in detecting observations that are randomly mixed with model generated samples. We report the test takers' accuracy rate in five tests.

:

S3.2 Power spectral analysis

We inspect the spatial structure of precipitation estimates by computing their average spectrum power as function of spatial frequencies and orientations. The computation steps are as follows:

Step 1: Transform the precipitation field data to frequency domain, using Fast Fourier Transform.

Step 2: Compute the power spectrum by taking the squared magnitude of the Fourier Transform coefficients.

Next, for radial averaged power spectrum analysis:

Step 1: Define a set of concentric circles centered at the origin of the frequency domain, along each radial line, calculate the average power by averaging the power spectrum values corresponding to the points intersected by the line.

Step 2: Plot the average power values against the corresponding radial frequency.

For orientation averaged power spectrum analysis:

Step 1: Define a set of concentric circles centered at the origin of the frequency domain, along each orientation angle, calculate the average power by averaging the power spectrum values corresponding to the points through this orientation angle.

Step 2: Plot the average power values against the corresponding orientation angle.

S3.3 Spread-Skill Correlation (SSC)

The spatial correlation coefficient between the standard deviation of model's ensemble, and the mean absolute error of model's ensemble mean:

:

$$SSC = \frac{\sum_{i=1}^{n} (\sigma_i - \bar{\sigma})(\epsilon_i - \bar{\epsilon})}{\sqrt{\sum_{i=1}^{n} (\sigma_i - \bar{\sigma})^2 \cdot \sum_{i=1}^{n} (\epsilon_i - \bar{\epsilon})^2}}$$
(19)

where σ_i is the standard deviation of model's ensemble at grid *i*:

$$\sigma_i = \sqrt{\frac{1}{J-1} \sum_{j=1}^{J} (\hat{y}_i^j - \bar{\hat{y}}_i)^2}$$
(20)

 ϵ_i is the mean absolute error of model's ensemble mean at grid *i*:

$$\sigma_i = |\bar{y}_i - y_i| \tag{21}$$

J is ensemble size; \hat{y}_i^j is the *j*th ensemble estimate at grid *i*; $\bar{\hat{y}}_i$ is ensemble mean estimate; y_i is observation.

S3.4 Coverage Ratio (CR)

The percentage that grid observation falls into the coverage of ensemble spread.

S3.5 Pearson correlation coefficient (r)

The Pearson correlation coefficient (r) between ensemble mean prediction \hat{y} and observation y is calculated as follows:

$$r = \frac{\sum_{i=1}^{n} (\hat{y}_i - \bar{y})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^{n} (\hat{y}_i - \bar{y})^2 \cdot \sum_{i=1}^{n} (y_i - \bar{y})^2}}$$
(22)

S3.6 Root mean squared error (RMSE)

The root mean square error (RMSE) between ensemble mean prediction \hat{y} and observation y is calculated as follows:

$$RMSE = \sqrt{\frac{\sum_{i=1}^{n} (y_i - \hat{y}_i)^2}{n}}$$
(23)

S3.7 Continuous ranked probabilistic score (CRPS)

November 23, 2023, 5:12am

The continuous ranked probability score (CRPS) is defined as:

$$\operatorname{CRPS}(F, x) = \int_{-\infty}^{\infty} \left[F(\hat{y}) - \mathbb{I}(\hat{y} \ge y) \right]^2 \, dy \tag{24}$$

where $F(\hat{y})$ is the cumulative distribution function (CDF) of the predictive distribution, y is the observed value, and $\mathbb{I}(\cdot)$ is the indicator function.

:

References

- Ho, J., Jain, A., & Abbeel, P. (2020). Denoising diffusion probabilistic models. Advances in neural information processing systems, 33, 6840–6851.
- Ho, J., & Salimans, T. (2022). Classifier-free diffusion guidance. *arXiv preprint* arXiv:2207.12598.
- Kingma, D., Salimans, T., Poole, B., & Ho, J. (2021). Variational diffusion models. Advances in neural information processing systems, 34, 21696–21707.
- Luo, C. (2022). Understanding diffusion models: A unified perspective. arXiv preprint arXiv:2208.11970.
- Nichol, A. Q., & Dhariwal, P. (2021). Improved denoising diffusion probabilistic models. In *International conference on machine learning* (pp. 8162–8171).
- Skamarock, W. C., Klemp, J. B., Dudhia, J., Gill, D. O., Liu, Z., Berner, J., ... others (2019). A description of the advanced research wrf version 4. NCAR tech. note ncar/tn-556+ str, 145.
- Sohl-Dickstein, J., Weiss, E., Maheswaranathan, N., & Ganguli, S. (2015). Deep unsupervised learning using nonequilibrium thermodynamics. In *International conference* on machine learning (pp. 2256–2265).
- Song, Y., Sohl-Dickstein, J., Kingma, D. P., Kumar, A., Ermon, S., & Poole, B. (2020). Score-based generative modeling through stochastic differential equations. arXiv preprint arXiv:2011.13456.



:

Figure S1. Model architecture of the unconditional score estimating neural network. We embed the time information, and stack the time embedding as extra channel to all UNet blocks. Each contract block consists of a long chain of $\{C_{3\times3} + N + \text{ReLU}\}_3$, and a short chain of $\{C_{1\times1}\}_1$, concatenated as a residual block. $C_{n\times n}$ is convolution layer, with kernel receptive field of size $n \times n$, N is group normalization. Each expand block consists of a long chain of $\{R_2 + C_{3\times3} + N + \text{ReLU}\}_3$, and a short chain of $\{R_2, C_{1\times1}\}_1$, concatenating as a residual block, R_n resize the data by n times using linear interpolation. We start with channel size of 128, and double/shrink the channel size by 2 along each contract/expand block.



Figure S2. Model architecture of the conditional score estimating neural network, similar to the unconditional score estimating neural network, but includes the conditioning information from a UNet precipitation estimation using dynamical field as input.



Figure S3. UNet architecture. The model takes into input of resolved dynamical field information and static elevation information, and outputs a deterministic precipitation field estimate. The dynamical field information is provided by a 9-hour (including 3 previous/current/future hours), $8^{\circ} \times 8^{\circ}$ circulation field data, with 19 channels representing 19 dynamical variables. This dynamical field information is first pre-processed through 3D convolution blocks (bottom), and concatenated with preprocessed elevation information, before feeding into a 2D UNet. The UNet applies a convolution based contracting path to capture precipitation relevant dynamical field information, and a symmetric transposed convolution based expanding path to gradually refine precipitation field estimates. Skip connections between symmetrical convolution and transposed convolution blocks are applied to force deeper neural network layers to learn meaningful representations that are not well captured by shallower layers.



:

Figure S4. CVAE architecture. The encoder approximates $p(\mathbf{z}|\mathbf{x}, \mathbf{y})$ as a variational Gaussian distribution; the decoder approximates $p(\mathbf{y}|\mathbf{x}, \mathbf{z})$ using the conditioning information \mathbf{x} and the learned latent vector \mathbf{z} . The dimension of \mathbf{z} is 32.



Figure S5. CGAN architecture. We replicate the generator network G for 10 times, each receiving a different \mathbf{z} and a same conditioning information \mathbf{x} to create an ensemble member. The generator network G is trained to fool the discriminator network D, while drawing the ensemble mean close to the realized observation. The discriminator network D is a binary classifier, optimized to differentiate between generated samples and true samples.

Physical process	Option
Cloud microphysics	Lin (Purdue)
Cumulus	Zhang and McFarlane
Radiation	Rapid Radiative Transfer Model
Boundary layer	Yonsei University (YSU) PBL scheme
Surface	Noah Land Surface Mode

Table S1. Physics options for WRF simulation.



Typhoon Lekima 0000 UTC 04 Aug 2019–0000 UTC 12 Aug 2019

Figure S6. WRF nested Domains (27km/9km/3km) for Typhoon Lekima simulation, from 0000 UTC 04 August 2019 to 0000 UTC 12 August 2019. Color denotes maximum precipitation rate (mm/3h) through the simulation period.