

# Learning to infer weather states using partial observations

Jie Chao<sup>1</sup>, Baoxiang Pan<sup>2</sup>, Quanliang Chen<sup>3</sup>, Shangshang Yang<sup>4</sup>, Jingnan Wang<sup>5</sup>, congyi nai<sup>6</sup>, Yue Zheng<sup>7</sup>, Xichen Li<sup>8</sup>, Huiling Yuan<sup>9</sup>, Xi Chen<sup>2</sup>, Bo Lu<sup>10</sup>, and Ziniu Xiao<sup>2</sup>

<sup>1</sup>Chengdu University of Information Technology

<sup>2</sup>Institute of Atmospheric Physics, Chinese Academy of Sciences

<sup>3</sup>Plateau Atmosphere and Environment Key Laboratory of Sichuan Province, College of Atmospheric Science, Chengdu University of Information Technology

<sup>4</sup>Key Laboratory of Mesoscale Severe Weather Ministry of Education/School of Atmospheric Sciences, Nanjing University

<sup>5</sup>College of Computer, National University of Defense Technology

<sup>6</sup>Institute of Geographic Sciences and Natural Resources Research, Chinese Academy of Sciences

<sup>7</sup>ClusterTech Limited, Hong Kong

<sup>8</sup>International Center for Climate and Environment Sciences, Institute of Atmospheric Physics, Chinese Academy of Sciences

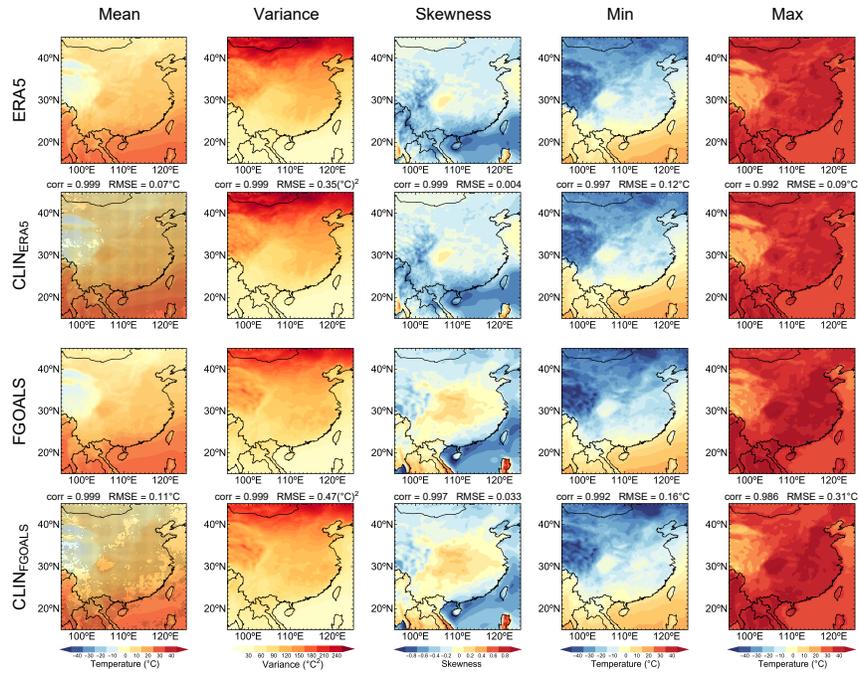
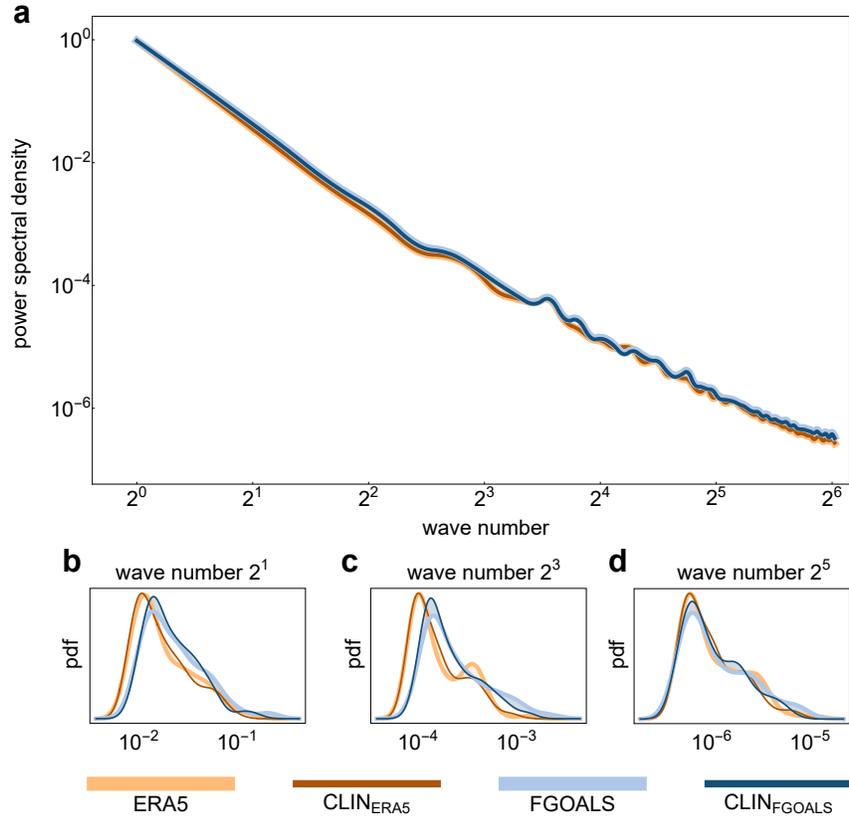
<sup>9</sup>Key Laboratory of Mesoscale Severe Weather/Ministry of Education, and School of Atmospheric Sciences, Nanjing University, Nanjing, China

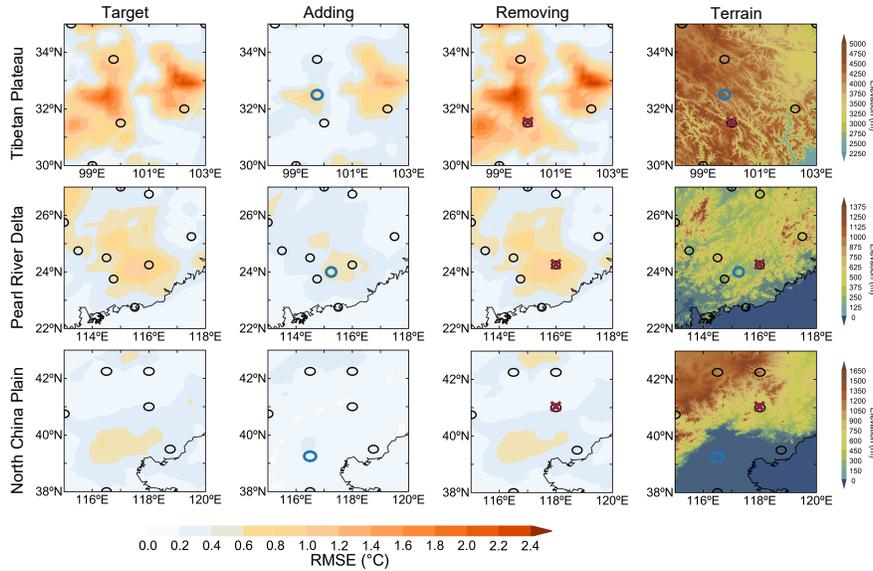
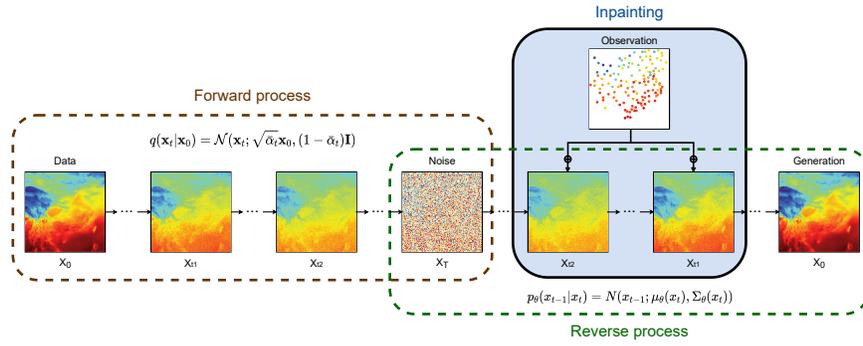
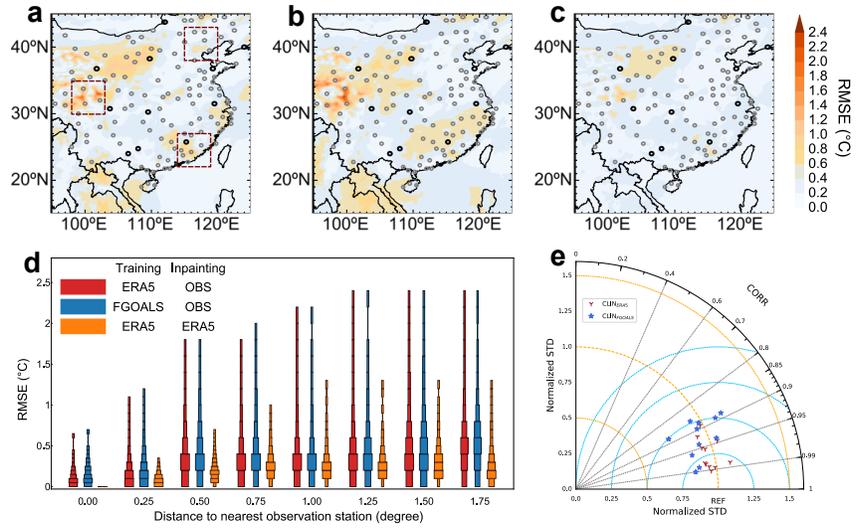
<sup>10</sup>National Climate Center, China Meteorological Administration

April 19, 2024

## Abstract

Accurate state estimation of the high-dimensional, chaotic Earth atmosphere marks a Sisyphean task, yet is indispensable for initiating weather forecast and gauging climate variability. While much effort is devoted to assimilating observations and forecasts to infer weather state, the inherent low-dimensional statistical structure in atmospheric circulation, shaped by geophysical laws and geographic boundaries, is underutilized as informative prior for state inference, or as reference for assessing representative of existing observations and planning new ones. We realize these potential by learning climatological distribution from climate reanalysis/simulation, using deep generative model. For a case study of estimating 2 m temperature spatial patterns, the learned distribution faithfully reproduces climatology statistics. A combination of the learned climatological prior with few station observations yields strong posterior of spatial pattern estimates, which are spatially coherent, faithful and adaptive to observation constraints, and uncertainty-aware. This allows us to evaluate each observation's value in reducing state estimation uncertainty, and guide optimal observation network design by pinpointing the most informative sites. Our study showcases how generative models can extract and utilize information produced in the chaotic evolution of climate system.





# Learning to infer weather states using partial observations

Jie Chao<sup>1,2</sup>, Baoxiang Pan<sup>2</sup>, Quanliang Chen<sup>1</sup>, Shangshang Yang<sup>2,3</sup>, Jingnan Wang<sup>2,4</sup>, Congyi Nai<sup>2,5</sup>, Yue Zheng<sup>6</sup>, Xichen Li<sup>2</sup>, Huiling Yuan<sup>3</sup>, Xi Chen<sup>2</sup>, Bo Lu<sup>7</sup>, Ziniu Xiao<sup>2</sup>

<sup>1</sup>School of Atmospheric Sciences, Chengdu University of Information Technology, Sichuan, China

<sup>2</sup>Institute of Atmospheric Physics, Chinese Academy of Science, Beijing, China

<sup>3</sup>Key Laboratory of Mesoscale Severe Weather, Ministry of Education, and School of Atmospheric Sciences, Nanjing University, Jiangsu, China

<sup>4</sup>College of Computer, National University of Defense Technology, Hunan, China

<sup>5</sup>Institute of Geographic Sciences and Natural Resources Research, Chinese Academy of Sciences, Beijing, China

<sup>6</sup>Clustertech LTD, Hong Kong, China

<sup>7</sup>National Climate Center, China Meteorological Administration, Beijing, China

## Key Points:

- Deep generative model enables accurate spatial interpolation of weather variables from sparse observations.
- The model generates probabilistic weather estimates with reliable uncertainty quantification by combining learned priors and observations.
- The model quantifies the value of observations for reducing uncertainty, guiding optimal observation network design.

---

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

78 multi-scale background information using forecasting models requires operational run-  
79 of large ensemble high-resolution numerical simulations, which is prohibitively expen-  
80 sive and burdensome (Toth et al., 2003; Palmer, 2017).

81 Is there extra information source for inferring the state of the high-dimensional,  
82 chaotic Earth atmosphere? It turns out that, the inherent low-dimensional statistical struc-  
83 ture in atmospheric circulation, shaped by the underlying geophysical laws and quasi-  
84 static geographic boundaries, can serve as an informative prior for state inference. The  
85 Earth climate system, like any other chaotic system, is an information producer: it grad-  
86 ually reveals the characteristic structure of its phase space at ever-ner scales (Gilpin,  
87 2024). By identifying and parameterizing this characteristic structure, we can potentially  
88 bypass the curse of high dimensionality, and make more efficient use of limited obser-  
89 vations for the state inference task.

90 Some pioneering works have explored this direction, leveraging the inherent struc-  
91 ture of climate data to fill in missing observations and rebuild historical climate records.  
92 For instance, Kadow et al. (2020) developed a partial convolution method to reconstruct  
93 historical global temperature patterns based on partial observations and climate simu-  
94 lation. Kannigier and Fiedler (2024) applied a similar methodology to restore the spa-  
95 tial extent of dust plumes in cloud-masked satellite images. Most of these practices con-  
96 sider deterministic models, which are designed for specific "reconstruction" problem con-  
97 figurations, yielding deterministic results regardless of whether observations can adequately  
98 constrain the estimation uncertainty. As a result, these methodologies generalize poorly  
99 to state inference tasks where the number or layout of observations change, fail to re-  
100 produce extremes or apply for scenarios where only limited observations are available.

101 A solution to these dilemmas is to shift from deterministic model to probabilistic  
102 model (B. Pan et al., 2021). Specifically, we prefer to build a probabilistic model that  
103 explicitly represents the inherent statistical structure of the atmosphere as revealed by  
104 climate observations or simulations. Thereafter, we hope to effectively and efficiently com-  
105 bine the learned climatological prior with incomplete observations, so as to obtain strong  
106 posterior of spatial pattern estimates. This problem setup poses two stringent require-  
107 ments on the underlying probabilistic model. First, the model must faithfully approx-  
108 imate the high-dimensional climatological distribution as generated by the chaotic evo-  
109 lution of climate dynamics. Second, the model must enable flexible probabilistic infer-  
110 ence, allowing us to efficiently obtain posterior atmospheric state estimates given arbi-  
111 trary observational constraints.

112 To fulfill these requirements, we resort to generative machine learning, in partic-  
113 ular, probabilistic diffusion models (Sohl-Dickstein et al., 2015; Ho et al., 2020; Song, Sohl-  
114 Dickstein, et al., 2020; Kingma et al., 2021). Probabilistic diffusion models learn to ap-  
115 proximate complex, high-dimensional probability distributions in an iterative manner,  
116 achieving unprecedented fitting capacity and controlling flexibility (B. Pan et al., 2023;  
117 Nai et al., 2024). To demonstrate the idea, we consider a case example of inferring the  
118 spatial pattern of 2 m temperature based on sparse observations from operational me-  
119 teorology stations. We learn probabilistic diffusion models to approximate the climato-  
120 logical distribution of 2 m temperature spatial patterns from climate reanalysis or sim-  
121 ulation data. After carefully assessing the model's ability to reproduce climatology, we  
122 develop tools to "inpaint" arbitrary observation constraints into the sample generation  
123 process, yielding probabilistic 2 m temperature spatial pattern estimates. Finally, we ap-

71 Deficiencies in observation render it an ill-posed task to estimate the state of the  
72 high-dimensional Earth atmosphere, calling for strong prior to achieve feasible solution.  
73 Forecasts from previous time steps are frequently applied to serve this mission, carry-  
74 ing information from previous step observations to the current step via a process-based  
75 model (Wang et al., 2000). As a result, the state estimation accuracy depends on an in-  
76 tricate interplay among model biases, background uncertainty, and observation error, which  
77 cannot be effectively disentangled or controlled (Law et al., 2015). Moreover, to provide  
78 multi-scale background information using forecasting models requires operational run  
79 of large ensemble high-resolution numerical simulations, which is prohibitively expen-  
80 sive and burdensome (Toth et al., 2003; Palmer, 2017).

81 Is there extra information source for inferring the state of the high-dimensional,  
82 chaotic Earth atmosphere? It turns out that, the inherent low-dimensional statistical struc-  
83 ture in atmospheric circulation, shaped by the underlying geophysical laws and quasi-  
84 static geographic boundaries, can serve as an informative prior for state inference. The  
85 Earth climate system, like any other chaotic system, is an information producer: it grad-  
86 ually reveals the characteristic structure of its phase space at ever-finer scales (Gilpin,  
87 2024). By identifying and parameterizing this characteristic structure, we can potentially  
88 bypass the curse of high dimensionality, and make more efficient use of limited obser-  
89 vations for the state inference task.

90 Some pioneering works have explored this direction, leveraging the inherent struc-  
91 ture of climate data to fill in missing observations and rebuild historical climate records.  
92 For instance, Kadow et al. (2020) developed a partial convolution method to reconstruct  
93 historical global temperature patterns based on partial observations and climate simu-  
94 lation. Kanngießer and Fiedler (2024) applied a similar methodology to restore the spa-  
95 tial extent of dust plumes in cloud-masked satellite images. Most of these practices con-  
96 sider deterministic models, which are designed for specific “reconstruction” problem con-  
97 figurations, yielding deterministic results regardless of whether observations can adequately  
98 constrain the estimation uncertainty. As a result, these methodologies generalize poorly  
99 to state inference tasks where the number or layout of observations change, fail to re-  
100 produce extremes or apply for scenarios where only limited observations are available.

101 A solution to these dilemmas is to shift from deterministic model to probabilistic  
102 model (B. Pan et al., 2021). Specifically, we prefer to build a probabilistic model that  
103 explicitly represents the inherent statistical structure of the atmosphere as revealed by  
104 climate observations or simulations. Thereafter, we hope to effectively and efficiently com-  
105 bine the learned climatological prior with incomplete observations, so as to obtain strong  
106 posterior of spatial pattern estimates. This problem setup poses two stringent require-  
107 ments on the underlying probabilistic model. First, the model must faithfully approx-  
108 imate the high-dimensional climatological distribution as generated by the chaotic evo-  
109 lution of climate dynamics. Second, the model must enable flexible probabilistic infer-  
110 ence, allowing us to efficiently obtain posterior atmospheric state estimates given arbi-  
111 trary observational constraints.

112 To fulfill these requirements, we resort to generative machine learning, in partic-  
113 ular, probabilistic diffusion models (Sohl-Dickstein et al., 2015; Ho et al., 2020; Song, Sohl-  
114 Dickstein, et al., 2020; Kingma et al., 2021). Probabilistic diffusion models learn to ap-  
115 proximate complex, high-dimensional probability distributions in an iterative manner,  
116 achieving unprecedented fitting capacity and controlling flexibility (B. Pan et al., 2023;  
117 Nai et al., 2024). To demonstrate the idea, we consider a case example of inferring the  
118 spatial pattern of 2 m temperature based on sparse observations from operational me-  
119 teorology stations. We learn probabilistic diffusion models to approximate the climato-  
120 logical distribution of 2 m temperature spatial patterns from climate reanalysis or simu-  
121 lation data. After carefully assessing the model’s ability to reproduce climatology, we  
122 develop tools to “inpaint” arbitrary observation constraints into the sample generation  
123 process, yielding probabilistic 2 m temperature spatial pattern estimates. Finally, we ap-

124 ply this methodology to evaluate each observation’s value in reducing state estimation  
 125 uncertainty, and guide optimal observation network design by pinpointing the most in-  
 126 formative sites.

## 127 2 Methodology

### 128 2.1 Data and problem setup

129 We consider the task of inferring the spatial pattern of 2 m temperature over East  
 130 Asia ( $15^\circ\text{N} - 45^\circ\text{N}, 95^\circ\text{E} - 125^\circ\text{E}$ ), using station observations covering  $\sim 1\%$  grids of  
 131 the considered region. To achieve this, we learn climatological distribution of 2 m tem-  
 132 perature spatial pattern using climate reanalysis or simulation data. The reanalysis data  
 133 are hourly,  $0.25^\circ$  2 m temperature data from the fifth-generation global climate and weather  
 134 reanalysis (ERA5) developed at European Centre for Medium-Range Weather Forecasts  
 135 (Hersbach et al., 2020, ECMWF). The simulation data are 3-hourly,  $0.25^\circ$  2 m temper-  
 136 ature historical simulation from the Flexible Global Ocean-Atmosphere-Land System Model  
 137 version f3-H (Bao et al., 2020, FGOALS-f3-H), which participates in the sixth phase of  
 138 the Coupled Model Intercomparison Project (Eyring et al., 2016, CMIP6). The station  
 139 observation data are obtained from the Chinese National Climatic Data Center (X. Pan  
 140 et al., 2021).

Formally, we denote the spatial pattern of 2 m temperature for the target region  
 as  $\mathbf{x}$ , which is a  $120 \times 120$  dimensional random variable here. Our objective is to ap-  
 proximate the distribution of  $\mathbf{x}$ , based on large number of samples from climate reanal-  
 ysis or simulation:

$$p_{\theta^*} = \arg \max_{p_\theta} \sum \log p_\theta(\mathbf{x}) \quad (1)$$

141 Here  $p_\theta$  is parameterized probability density function approximator,  $\theta^*$  is the optimal  
 142 parameter, optimized by maximizing the overall likelihood of  $p_\theta$  assigned to the train-  
 143 ing samples.

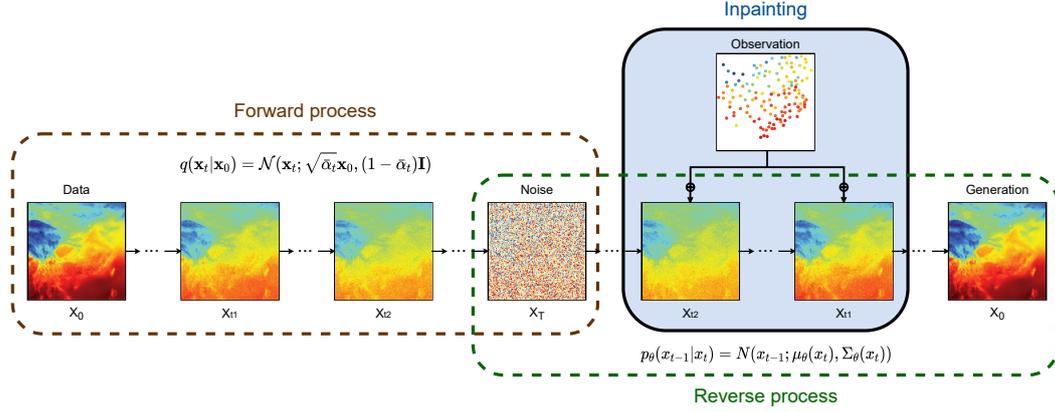
144 Given  $p_{\theta^*}$  and sparse observations, we need to provide probabilistic estimates of  
 145 2 m temperature spatial patterns, i.e.,  $p_{\theta^*}(\mathbf{x}|\mathbf{x} \odot \mathbf{m})$ . Here,  $\odot$  is dot product,  $\mathbf{m}$  is ob-  
 146 servation mask, with value 1/0 denoting the existence/absence of observations for each  
 147 geogrid.  $p_{\theta^*}(\mathbf{x}|\mathbf{x} \odot \mathbf{m})$  should yield samples that are spatially coherent and faithful to  
 148 observational constraints. Also,  $p_{\theta^*}(\mathbf{x}|\mathbf{x} \odot \mathbf{m})$  should offer accurate uncertainty quan-  
 149 tification. For instance, geogrids close to observation stations should typically have low  
 150 state estimate uncertainties, while distant ones have high uncertainties. Finally, we pre-  
 151 fer  $p_{\theta^*}(\mathbf{x}|\mathbf{x} \odot \mathbf{m})$  to be adaptive to changes in observation configurations, such as the  
 152 abortion or inclusion of observation stations, or rearrangement of station network lay-  
 153 out. Below we illustrate how to achieve these requirements using the proposed method-  
 154 ology.

### 155 2.2 Learning climatology with probabilistic diffusion model

156 We elucidate how to learn climatological distribution of the target random vari-  
 157 able using probabilistic diffusion model, thereafter leverage this learned prior for the in-  
 158 ference task (Sec. 2.3). For clarity, we only cover key steps necessary for establishing our  
 159 methodology. Details can be found in the literature referenced through the description.

160 To approximate a target distribution using probabilistic diffusion model, we train  
 161 a series of deep neural networks that can be chained to establish bijective mapping be-  
 162 tween the target distribution and a prior distribution (Sohl-Dickstein et al., 2015; Ho et  
 163 al., 2020). Specifically, we define the following Gaussian process:

$$q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{(1 - \beta_t)}\mathbf{x}_{t-1}, \beta_t\mathbf{I}) \quad (2)$$



**Figure 1.** Overview of the Climate Inpainting (CLIN) methodology. A pre-defined forward Gaussian process (left) turns distribution of target climate variable into a prior distribution, i.e., standard Gaussian. A learned reverse Gaussian process (right) turns the prior distribution into the distribution of the target climate variable. We “inpaint” sparse observations throughout the reverse Gaussian process (right top), so as to obtain spatial pattern estimates of the target variable.

164 Here  $p(\mathbf{x}_0) = p(\mathbf{x})$ , which is the target distribution;  $p(\mathbf{x}_T)$  is the prior distribution; we  
 165 bridge  $\mathbf{x}_0$  and  $\mathbf{x}_T$  using  $\mathbf{x}_{t \in [1, T]}$ , which are latent variables with increasing noise level;  
 166  $\mathcal{N}$  is Gaussian distribution;  $\mathbf{I}$  is identity matrix;  $\beta_t$  is diffusion coefficient, which is pre-  
 167 defined so that, give large enough  $T$ ,  $p(\mathbf{x}_T | \mathbf{x}_0)$  is drawn close to  $p(\mathbf{x}_T)$ , which is  $\mathbf{x}_0$  ag-  
 168 nostic. This setup offers analytical solution for  $p(\mathbf{x}_{t+\tau} | \mathbf{x}_t), \forall \tau \in [0, T - t], t \in [0, T]$ ,  
 169 facilitating convenient inference as detailed in Sec. 2.3.

To achieve generative modeling, we reverse Eq. 2 using the following variation distributions:

$$p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \mu_\theta, \Sigma_\theta) \quad (3)$$

170 Here  $\Sigma_\theta$  is represented as an interpolation between its analytical lower and upper bound  
 171 (Dhariwal & Nichol, 2021);  $\mu_\theta$  can be optimized by maximizing the variational lower bound  
 172 (ELBO) on the log-likelihood of the training samples (Sohl-Dickstein et al., 2015; Kingma  
 173 et al., 2021). In practice, we represent  $\mu_\theta$  as function of neural network parameteriza-  
 174 tion for  $\nabla p(\mathbf{x}_t | \mathbf{x}_0)$ , which is known as the *score function* (Song, Garg, et al., 2020; Song,  
 175 Sohl-Dickstein, et al., 2020). This simplifies the ELBO objective function to the follow-  
 176 ing form:

$$L = \mathbb{E}_{t \in [1, T], \mathbf{x}_0 \sim p(\mathbf{x}_0)} \|\nabla p(\mathbf{x}_t | \mathbf{x}_0) - \epsilon_\theta\|^2 \quad (4)$$

177 Here  $\epsilon_\theta$  is a neural network parameterization for  $\nabla p(\mathbf{x}_t | \mathbf{x}_0)$ . Given the trained score es-  
 178 timates, we can derive  $p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \mu_\theta, \Sigma_\theta)$  and sample it, starting with  $p(\mathbf{x}_T)$ ,  
 179 ending with  $p(\mathbf{x}_0)$ .

### 2.3 CLIN: inferring weather states using partial observations

180  
 181 We combine the learned climatology prior with station observations to infer the pos-  
 182 terior probability distribution of the target variable, using a *repainting* methodology (Lugmayr  
 183 et al., 2022; Zhang et al., 2023). Specifically, given a pre-trained diffusion model that se-  
 184 quentially applies  $p_{\theta^*}(\mathbf{x}_t | \mathbf{x}_{t+1}) = \mathcal{N}(\mathbf{x}_t; \mu_{\theta^*}, \Sigma_{\theta^*})$  to transform  $p(\mathbf{x}_T)$  to  $p(\mathbf{x}_0)$ , within  
 185 a pre-selected time window of  $\Omega$ , for grid points where we have observations, we replace

186 values of  $\mathbf{x}_t$  with observations noisified to time step  $t$ , by sampling  $p(\mathbf{x}_t \odot \mathbf{m} | \mathbf{x}_0 \odot \mathbf{m})$ .  
 187 This replacement does not consider the generated parts of  $\mathbf{x}_t$ , therefore, the observations  
 188 could not explicitly constrain the variability of unobserved parts.

189 To address this issue, for any  $t \in \Omega$ , after the replacement, instead of progress-  
 190 ing to  $t - 1$  directly, we rewind to time step  $t - \tau$  by sampling  $p(\mathbf{x}_{t-\tau} | \mathbf{x}_t)$ . We there-  
 191 after repeat the denoising steps from  $t - \tau$  to  $t$  for  $k$  rounds, and carry out observation  
 192 replacement for  $\mathbf{x}_t$  at each round. This allows us to jointly modify both observed and  
 193 unobserved regions throughout the denoising steps, yielding generated samples that are  
 194 spatially coherent, faithful and adaptive to observation constraints, and uncertainty-aware.  
 195 This methodology is referred to as *inpainting*, we hence name our methodology as CLIN,  
 196 short for Climate Inpainting. A formal algorithm description is given below. Details for  
 197 data processing, neural network architecture, hyperparameters for training and inference,  
 198 are given in Supporting Information.

---

**Algorithm 1** CLIN

---

**Require:** trained diffusion model  $p_{\theta^*}$ , observations  $\mathbf{x}_0 \odot \mathbf{m}$ , repainting time step set  $\Omega$ ,  
 rewinding step  $\tau$ , rewinding round  $K$   
**Ensure:** observation constrained, spatially coherent sample  $\mathbf{x}_0$

- 1: Initialize  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 2: **for**  $t = T - 1, \dots, 1$  **do**
- 3:      $\mathbf{x}_t \sim p_{\theta^*}(\mathbf{x}_t | \mathbf{x}_{t+1})$  ▷ Reverse sampling
- 4:     **if**  $t \in \Omega$  **then:**
- 5:         **for**  $k = 1, \dots, K$  **do**
- 6:              $\mathbf{x}_t^{\text{obser}} \sim p(\mathbf{x}_t \odot \mathbf{m} | \mathbf{x}_0 \odot \mathbf{m})$
- 7:              $\mathbf{x}_t \leftarrow \mathbf{x}_t \odot (\mathbf{I} - \mathbf{m}) + \mathbf{x}_t^{\text{obser}}$  ▷ Condition on observations
- 8:              $\mathbf{x}_{t+\tau} \sim p(\mathbf{x}_{t+\tau} | \mathbf{x}_t)$  ▷ Rewind in time by  $\tau$  steps
- 9:             **for**  $i = t + \tau - 1, \dots, t$  **do**
- 10:                  $\mathbf{x}_i \sim p_{\theta^*}(\mathbf{x}_i | \mathbf{x}_{i+1})$  ▷ Reverse sampling within a rewinding round
- 11:             **end for**
- 12:         **end for**
- 13:     **end if**
- 14: **end for**
- 15: **return**  $\mathbf{x}_0$

---

199 **3 Results**

200 The accuracy for state estimation depends on 1) how well we can approximate the  
 201 climatological distribution, and 2) based on a learned climatological prior, how well we  
 202 can combine it with limited observations to obtain probabilistic state estimates. Below  
 203 we assess model’s performance for these two aspects (Sec. 3.1 and 3.2). We further em-  
 204 ploy the model to quantify the extent to which observations reduce uncertainty in state  
 205 estimation, offering insights for optimal observation design (Sec. 3.3).

206 **3.1 Climatology**

207 We compare grid-scale and field-scale statistics of 10,000 reference/generated sam-  
 208 ples to evaluate how well the probabilistic diffusion models reproduce their training data’s  
 209 climatology. Two models trained with climate reanalysis (ERA5) and historical climate  
 210 simulation (FGOALS) data, hereafter referred to as  $\text{CLIN}_{\text{ERA5}}$  and  $\text{CLIN}_{\text{FGOALS}}$ , are  
 211 deployed and evaluated.

212 The grid-scale assessment considers the mean, variance, skewness, minimum, and  
 213 maximum of climatological distribution at each grid (Fig. 2). These statistics from ERA5  
 214 (Fig. 2 Row 1) and FGOALS (Fig. 2 Row 3) generally agree well, due to shared constraints  
 215 from geophysical laws and geographic boundaries. The key spatial patterns are the lat-  
 216 itudinal gradient, the influence of topography (e.g., the Tibetan Plateau), and the land-  
 217 sea contrast, which are most evident in the mean, minimum and maximum maps. The  
 218 variance and skewness maps reveal more regional variations. A notable discrepancy is  
 219 that, compared to ERA5, FGOALS tends to hold larger skewness for most of the land  
 220 regions in Southern China and Philippine Island, implying a more frequent present of  
 221 high 2 m temperature for these regions.

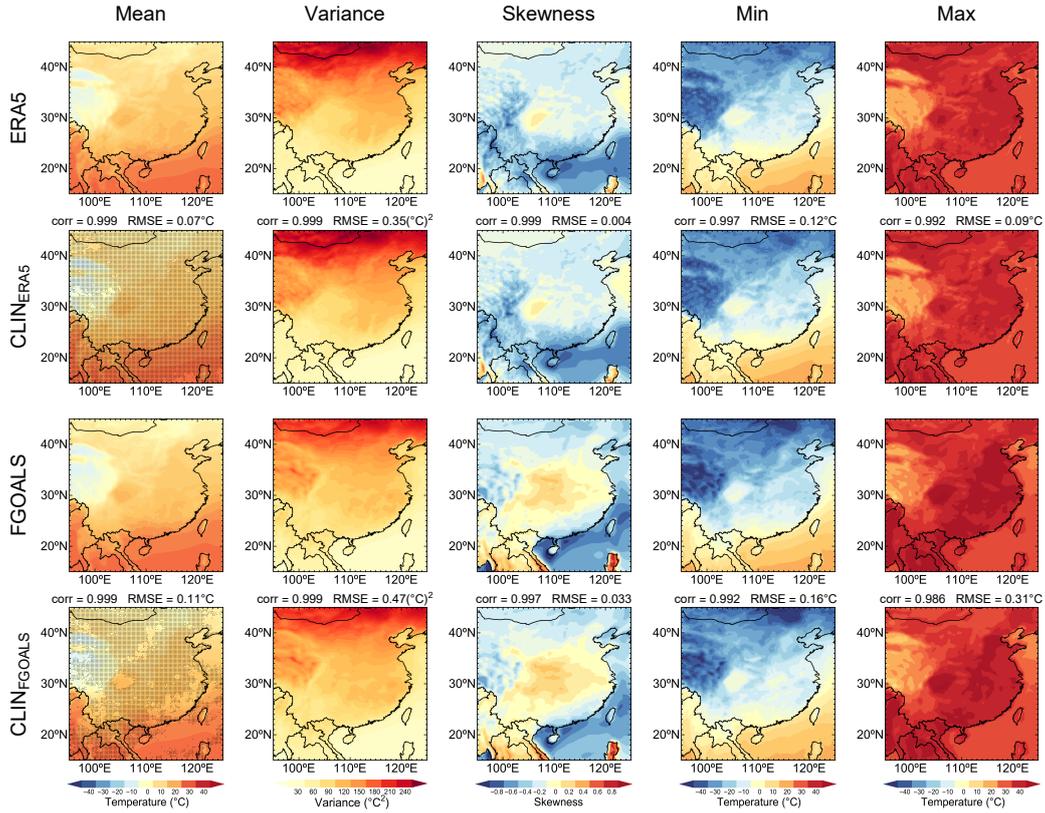
222  $CLIN_{ERA5}$  (Fig. 2 Row 2) and  $CLIN_{FGOALS}$  (Fig. 2 Row 4) can well reproduce the  
 223 considered statistics of their training data, achieving high spatial correlation coefficient  
 224 ( $\sim 0.99$ ) and low root mean squared error ( $\sim 0.1^\circ\text{C}$ ) in matching these statistics. Be-  
 225 sides reproducing the large scale patterns, both models accurately capture high frequency  
 226 local variations influenced by complex topography, such as for mountainous regions and  
 227 coastal areas. Also, the climatology difference between ERA5 and FGOALS are well re-  
 228 produced by the corresponding CLIN models.

229 We further carry out grid-wise Kolmogorov-Smirnov tests to assess whether the gen-  
 230 erated and referential samples are likely to have come from the same underlying distri-  
 231 bution: 96/76% grid points (stippled grids in Fig. 2) within the considered region pass  
 232 a 95% confidence interval test for the  $CLIN_{ERA5}$  and  $CLIN_{FGOALS}$  model. These results  
 233 suggest that the CLIN model can well reproduce climatological distribution of its train-  
 234 ing data at grid scale.

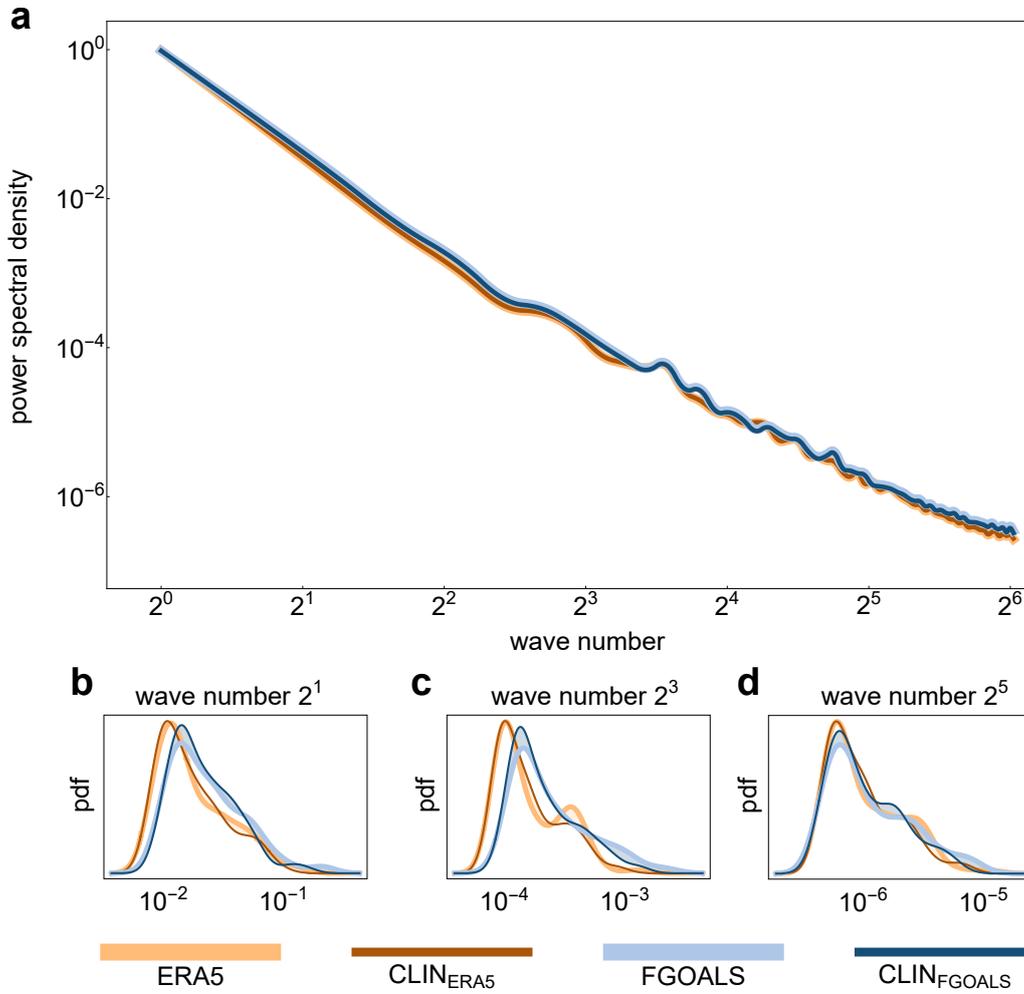
235 We hereafter compare the referential and generated distributions using field-scale  
 236 statistics. We first examine the linear spatial structure of the 2 m temperature spatial  
 237 patterns using a principal component analysis (Supporting Information Fig. S2): we de-  
 238 compose the spatial pattern of the target random variable into a set of orthogonal modes  
 239 that capture the maximum amount of variance, and compare the spatial modes (Em-  
 240 pirical Orthogonal Functions, EOFs), as well as the variance explained by these modes.  
 241 For ERA5, the first to third leading principal components explained 90/2.7/2.0% of the  
 242 total variance. While for  $CLIN_{ERA5}$ , the first to third leading principal components ex-  
 243 plained 91/2.6/1.5% of the total variance, which closely matches results for the ERA5  
 244 referential data. More importantly, we obtain spatial correlation coefficient of 0.994/0.990/0.986  
 245 between the first to third EOF of ERA5 and  $CLIN_{ERA5}$ . While the spatial modes of FGOALS  
 246 differs considerably with ERA5,  $CLIN_{FGOALS}$  closely matches FGOALS: the first to third  
 247 leading principal components explained 83.6/5.1/2.2% or 83.9/4.9/2.1% of the total vari-  
 248 ance for FGOALS or  $CLIN_{FGOALS}$ . The spatial correlation coefficient between the first  
 249 to third EOF of FGOALS and  $CLIN_{FGOALS}$  are 0.999/0.997/0.994. These results sug-  
 250 gest that the CLIN model can well reproduce the linear spatial mode of the considered  
 251 climatological distribution.

252 Lastly, we examine the distribution of spatial variability across different spatial scales  
 253 in the referential/generated dataset: we carry out 2D Fourier transform on the referen-  
 254 tial/generated samples, and draw the radial averaged squared magnitude of the complex  
 255 Fourier coefficients as function of wave numbers (Fig. 3). The radially averaged power  
 256 spectrum density of the considered referential and generated data samples follow a sim-  
 257 ilar power-law scaling, suggesting that the CLIN model can well reproduce the spatial  
 258 variability across scales.

259 To sum up, the analysis of both grid-scale and field-scale statistics demonstrates  
 260 that the CLIN methodology accurately reproduces the essential characteristics and pat-  
 261 terns of the climatological distribution present in the training data. We can thereafter  
 262 leverage this learned climatological prior for the state inference task.



**Figure 2.** Grid-scale comparison of climatological statistics for climate reanalysis (ERA5, Row 1), climate simulation (FGOALS, Row 3), and probabilistic diffusion models trained using these datasets (CLIN<sub>ERA5</sub>, Row 2; and CLIN<sub>FGOALS</sub>, Row 4). The considered statistics are mean, variance, skewness, minimum, and maximum. The spatial correlation coefficient (corr) and root mean squared error (RMSE) between the referential dataset statistics and generated dataset statistics are labeled. Stipples denote grids that pass the Kolmogorov-Smirnov test at 95% confidence interval.



**Figure 3.** Radial averaged power spectrum density as function of wave number for 2 m temperature spatial pattern. **a:** results for ERA5, FGOALS, CLIN<sub>ERA5</sub>, and CLIN<sub>FGOALS</sub> averaged over 100 ensemble members. **b-d:** probability distribution of power spectrum density at wave number 2<sup>1</sup>, 2<sup>3</sup>, 2<sup>5</sup> for ERA5, FGOALS, CLIN<sub>ERA5</sub>, and CLIN<sub>FGOALS</sub>.

263

### 3.2 Inferring weather states using partial observations

264

265

266

267

268

269

270

271

Given a learned climatological prior, we assess how well we can combine it with partial observations to obtain probabilistic estimate of the 2 m temperature spatial patterns. The climatological priors are probabilistic diffusion models trained using climate reanalysis (ERA5) and climate simulation (FGOALS) data. The observations are from 131 operational meteorological stations across China. We randomly select 120 of these stations to inpaint into the generation process, and leave the rest 11 stations for test. For regions without station observations, we consider ERA5 data as benchmark. Below we report case example results (Sec. 3.2.1) and a 1-year round skill assessment (Sec. 3.2.2) .

272

#### 3.2.1 Case study

273

274

275

276

277

278

279

280

We consider four case examples covering different hours of a day and different seasons (Fig. 4). To make probabilistic inference of spatial patterns using partial observations, we gradually inpaint station observations into the generation process of  $CLIN_{ERA5}$  and  $CLIN_{FGOALS}$ , creating 100 ensemble members for each model and each case. We report the ERA5 spatial pattern (Fig. 4 Row 1), the ensemble mean (Fig. 4 Row 2 and 5), the standard deviation of the ensemble (Fig. 4 Row 3 and 6), the mean squared error between ERA5 and the ensemble members (Fig. 4 Row 4 and 7) for  $CLIN_{ERA5}$  and  $CLIN_{FGOALS}$ .

281

282

283

284

285

286

Both the repainted  $CLIN_{ERA5}$  and  $CLIN_{FGOALS}$  ensemble mean results closely match the ERA5 spatial pattern, regarding latitudinal gradient, influence of topography, and the land-sea contrast, yielding spatial correlation coefficient of  $0.980 \pm 0.02 / 0.977 \pm 0.02$  for the four considered case examples. These results suggest that the proposed methodology allows effectively propagation of information from limited ( $\sim 1\%$ ) observed locations to a broad range of unobserved parts.

287

288

289

290

291

292

293

294

295

296

297

298

299

300

Next, we test if the CLIN methodology offers reliable uncertainty quantification (Fig. 4 Row 3 and 6). A larger ensemble variance indicates greater uncertainty in the estimate, while a smaller variance suggests more confidence in the estimate. As is expected, geogrids close to observation stations tend to have low ensemble variance, while distant ones may have relatively higher ensemble variance. The information constraint from observations may be blocked by topography, such as for Tibetan Plateau and Tian Shan Mountains. While for plain regions, we can expect a larger extension of observation constraints. We further examine the relationship between the spread of the ensemble members and their estimation skill, by computing the correlation between ensemble variance and ensembles' mean squared error score. The high spread skill correlation for  $CLIN_{ERA5}$  ( $0.90 \pm 0.08$ ) and  $CLIN_{FGOALS}$  ( $0.94 \pm 0.04$ ) suggest that ensemble spread is a good predictor of model's estimation skill. This means that the CLIN model can capture the underlying uncertainties and provide reliable estimates of spatial estimation confidence.

301

302

303

304

305

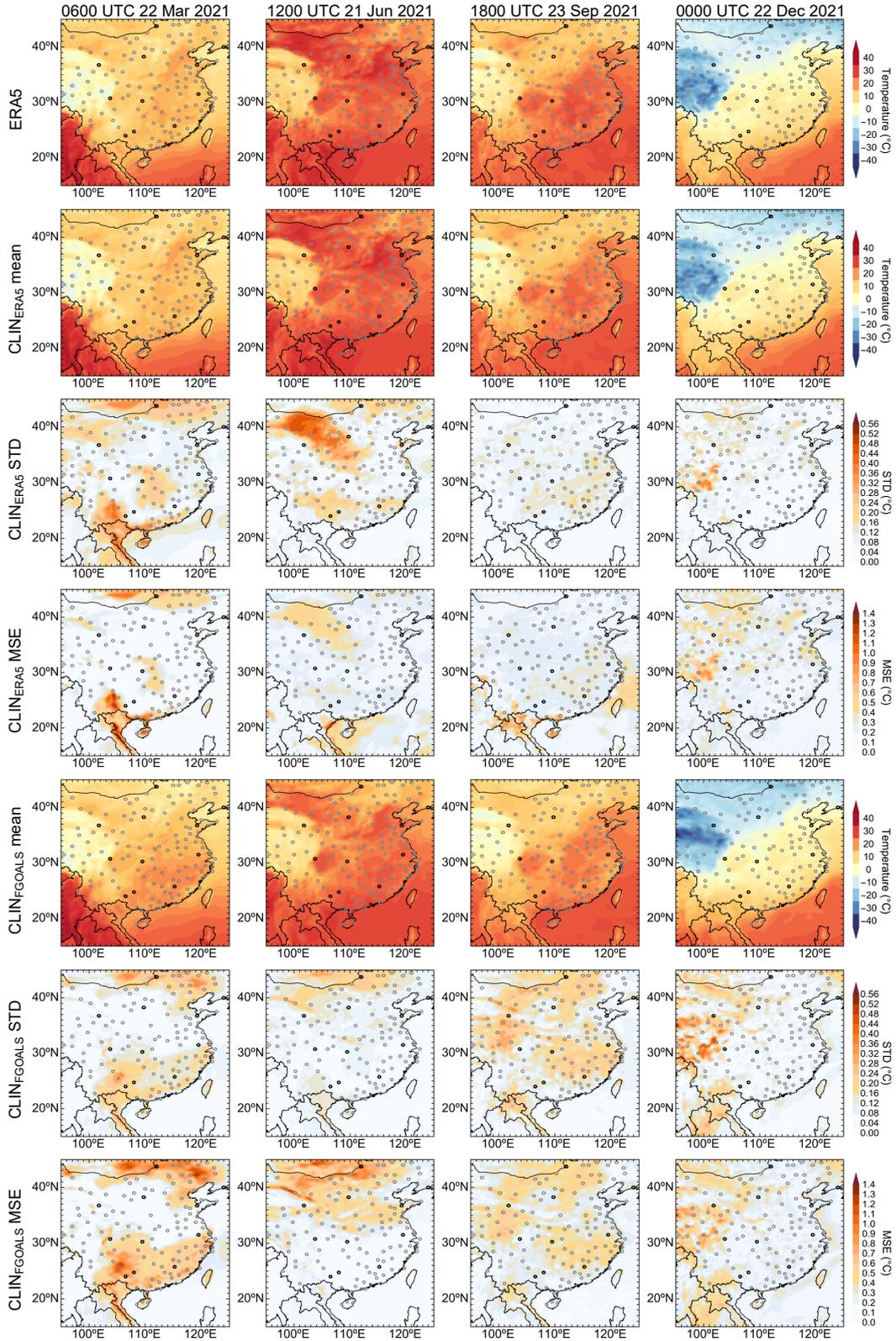
306

307

308

309

To sum up, the case studies confirm that the CLIN methodology can make successful probabilistic inference of 2 m temperature spatial patterns using limited observations. The results are spatially coherent, well-constrained by observations, and offer reliable uncertainty quantification. It is worth noting that there are unneglectable mismatches between station observations and ERA5/FGOALS, regarding either climatological statistics or values. These mismatches introduce domain shift error, which is frequently encountered as we deploy a machine learning model in real-world scenarios where the data distribution differs from the training data. Below we dissect this error source by inpainting with different data sources in a 1-year round evaluation.



**Figure 4.** Case examples for probabilistic inference for 2 m temperature spatial pattern using partial observations. For  $CLIN_{ERA5}$  and  $CLIN_{FGOALS}$ , 100 ensemble members are created by repainting observations. The ERA5 spatial pattern (Row 1), the ensemble mean (Row 2 and 5), the standard deviation of the ensemble (Row 3 and 6), the mean squared error between ERA5 and the ensemble members (Row 4 and 7) for  $CLIN_{ERA5}$  and  $CLIN_{FGOALS}$  are plotted.

310

### 3.2.2 Skill evaluation

311

312

313

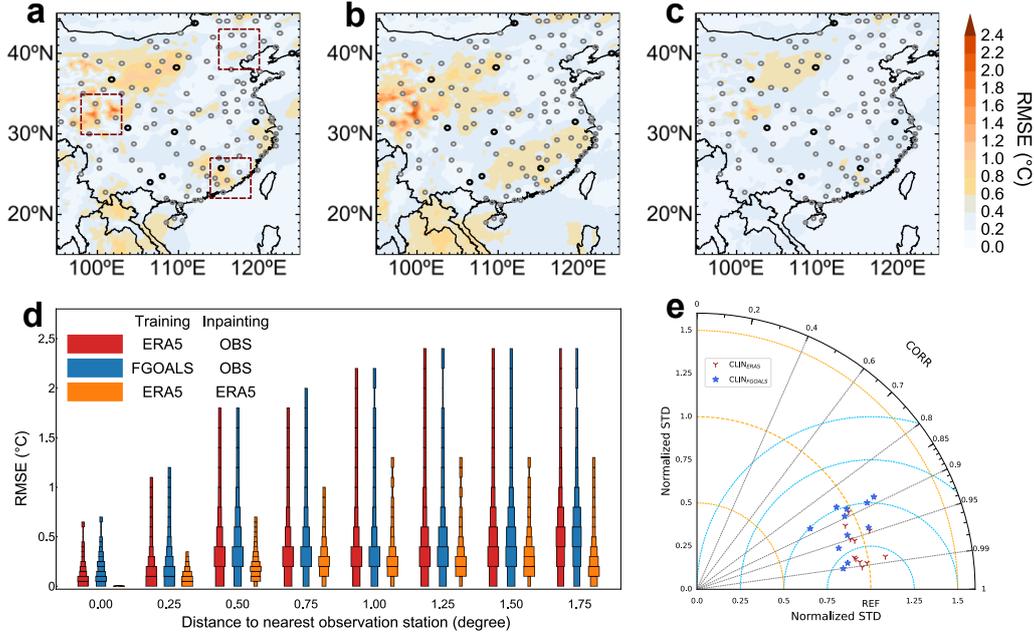
314

315

316

317

We conduct a year-long evaluation of the models' performance in inferring spatial patterns, using data from Year 2021, which are not included in the models' training process. We compare ERA5 with  $CLIN_{ERA5}$  and  $CLIN_{FGOALS}$ , both inpainted using station observations, and present the spatial distribution of their RMSE in Fig. 5a and Fig. 5b. To further investigate different uncertainty sources in the state inference task, we also consider inpainting  $CLIN_{ERA5}$  using ERA5 data at the observation stations. The RMSE between this inpainted  $CLIN_{ERA5}$  and the ERA5 whole-field data is shown in Fig. 5c.



**Figure 5.** Skill evaluation for CLIN models to estimate spatial pattern of 2 m temperature using data for Year 2021. **a:** root mean squared error (RMSE) between ERA5 reanalysis and  $CLIN_{ERA5}$  inpainted using station observations; **b:** RMSE between ERA5 reanalysis and  $CLIN_{FGOALS}$  inpainted using station observations; **c:** RMSE between ERA5 reanalysis and  $CLIN_{ERA5}$  inpainted using ERA5 data at station observations; **d:** distribution of RMSE as function of grid's distance to nearest observation station for the three considered methods; **e:** Taylor diagram comparing the left-out station observations with  $CLIN_{ERA5}$  (orange) and  $CLIN_{FGOALS}$  (blue) results. Both  $CLIN_{ERA5}$  and  $CLIN_{FGOALS}$  are constrained by 120 station observations here. We delineate three representative regions to evaluate the value of observations in Sec. 3.3

318

319

320

321

322

323

324

325

326

The RMSE between ERA5 and observation inpainted  $CLIN_{ERA5}/CLIN_{FGOALS}$  is  $0.25 \pm 0.21^\circ\text{C}/0.31 \pm 0.20^\circ\text{C}$ , suggesting that the CLIN methodology enables accurate spatial pattern estimates. Both models exhibit low uncertainty in plain terrain regions or over the ocean, despite that no ocean observations were applied. This suggests that the learned climatological prior effectively captures the spatial patterns and variability in these regions, allowing the models to make confident estimates using limited and far-away observational constraints. On the other hand, both models exhibit higher uncertainty in regions with complex terrain, such as the Tibetan Plateau and the mountainous areas of Southeast China. Additionally, land areas with complicated terrain but lack-

327 ing observational constraints, such as Southeast Asia, also show large uncertainty in the  
 328 model estimates.

329 The uncertainty in state inference comes from the following three sources (Tab. 1).  
 330 The first is domain shift error, which is due to distribution mismatch among data ap-  
 331 plied for model training, data applied for inpainting, and data applied for skill evalua-  
 332 tion. The second is model error, which is due to the approximation/optimization/statistical  
 333 error in applying probabilistic diffusion model to fit climatological prior, or due to er-  
 334 rors in inpainting. These two types of uncertainties are *epistemic*, as they could be re-  
 335 duced by gathering more data, improving the model, or incorporating knowledge about  
 336 data distribution differences. The third source of uncertainty is intrinsic/aleatoric, which  
 337 is due to existence of multiple plausible spatial patterns given partial observational con-  
 338 straints, reflecting the inherent randomness in the system being modeled.

339 To disentangle these uncertainty sources, we consider the following comparisons.

- 340 1. We compare the RMSE of  $\text{CLIN}_{\text{ERA5}}$  (Fig. 5a) and  $\text{CLIN}_{\text{FGOALS}}$  (Fig. 5b).  $\text{CLIN}_{\text{ERA5}}$   
 341 achieves an overall lower RMSE, which can be attributed to a relieved domain shift  
 342 error from the following two aspects: a. compared to FGOALS, ERA5 better matches  
 343 the “true” climatology as partially revealed by the scattered observations; b. we  
 344 consider ERA5 data as “ground truth” for evaluating model performance, which  
 345 gives advantage to CLIN model trained using ERA5 data.
- 346 2. We compare  $\text{CLIN}_{\text{ERA5}}$  inpainted using observation data (Fig. 5a) and  $\text{CLIN}_{\text{ERA5}}$   
 347 inpainted using scattered ERA5 data (Fig. 5c). The latter achieve significantly  
 348 lower RMSE ( $0.19 \pm 0.12^\circ\text{C}$ ), suggesting a relatively low model error and a rel-  
 349 atively low intrinsic uncertainty of the considered task. The difference between  
 350 these cases highlights the domain shift error as the observation distribution dif-  
 351 fers from ERA5.
- 352 3. We compare the performance of  $\text{CLIN}_{\text{ERA5}}$  and  $\text{CLIN}_{\text{FGOALS}}$  in predicting the  
 353 observations at test stations that are excluded during repainting (Fig. 5e). For these  
 354 test stations, both  $\text{CLIN}_{\text{ERA5}}$  and  $\text{CLIN}_{\text{FGOALS}}$  results show high correlation co-  
 355 efficient (0.87-0.99) and low root mean squared error (0.2-0.4 $^\circ\text{C}$ ) with the obser-  
 356 vations, with  $\text{CLIN}_{\text{ERA5}}$  performing slightly better than  $\text{CLIN}_{\text{FGOALS}}$ ;  $\text{CLIN}_{\text{ERA5}}$   
 357 holds a normalized standard deviation close to 1, which closely matches the ob-  
 358 servations, while  $\text{CLIN}_{\text{FGOALS}}$  holds a normalized standard deviation slightly less  
 359 than 1, suggesting a smaller temporal variability.

**Table 1.** Uncertainty sources for state inference using partial observations

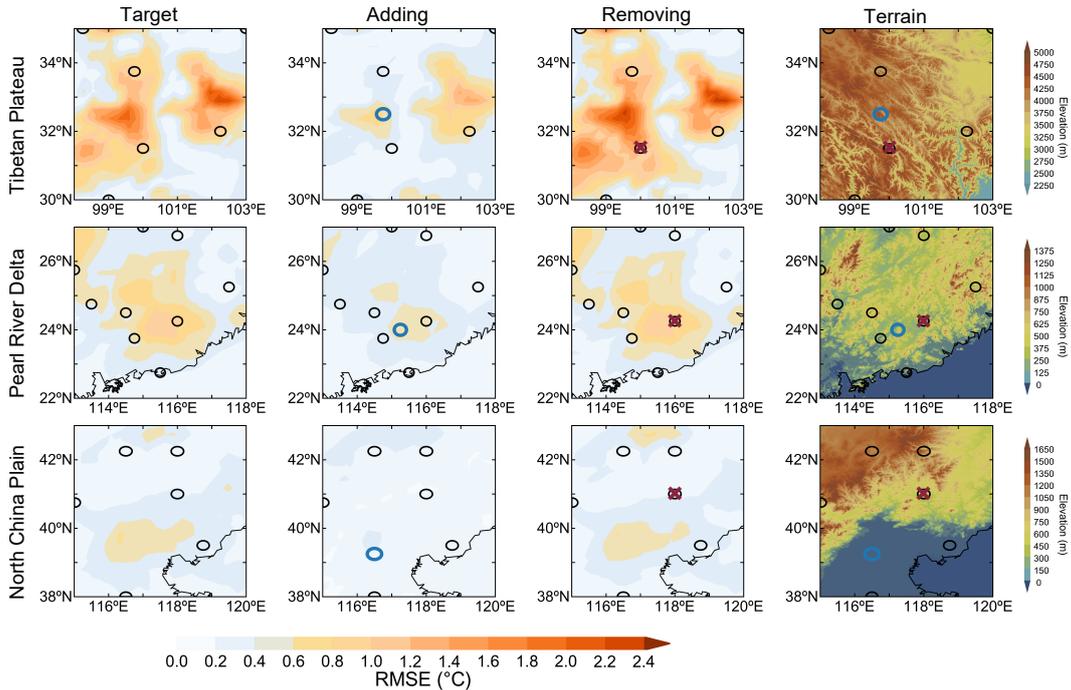
Uncertainty source	Type	Illustration
Domain shift	Epistemic	Distribution mismatch among data applied for model training, data applied for inpainting, and data applied for skill evaluation.
Model error	Epistemic	1. Approximation/optimization/statistical error in fitting climatological prior. 2. Error in constraining the prior with observations.
Intrinsic uncertainty	Aleatoric	Existence of multiple plausible spatial patterns given observational constraints.

360 Finally, we quantify the spatial extension of observational constraints by showing  
 361 models’ RMSE skill as function of grid’s distance to nearest observation station (Fig. 5d).  
 362 We consider  $\text{CLIN}_{\text{ERA5}}$  inpainted using observation data and ERA5 data, as well as  $\text{CLIN}_{\text{FGOALS}}$   
 363 inpainted using observation data. For all these cases, models’ performances at an arbi-  
 364 trary grid depends closely on the grid’s proximity to observations. Meanwhile, there is  
 365 large variation of models’ RMSE skills for grids that are at least  $1^\circ$  away from any ob-  
 366 servation stations. Below we further investigate the value of individual observations in

367 constraining the variability of its nearby spatial patterns, and offer guidelines for bet-  
 368 ter observation planning.

### 369 3.3 On the value of observations

370 We apply the CLIN methodology to quantify the value of observations in constrain-  
 371 ing state estimation uncertainty, using three representative regions delineated in Fig. 5a.  
 372 To achieve this, we add or remove observational stations and evaluate the impact on the  
 373 estimation error (Fig. 6). Here, the first column shows the RMSE spatial pattern for the  
 374 original  $\text{CLIN}_{\text{ERA5}}$  model estimates in each target region; the second column (Adding)  
 375 demonstrates the impact of adding an observation station in a high-error area; the third  
 376 column (Removing) illustrates the effect of removing an existing observation station; the  
 377 fourth column (Terrain) provides a topographical context for each target region.



**Figure 6.** Evaluation of CLIN in reconstructing 2m temperature spatial pattern using different observation setups. Column 1: RMSE between  $\text{CLIN}_{\text{ERA5}}$  inpainted using observation data and ERA5 for three selected regions delineated in Fig. 5. Column 2: RMSE after including a pseudo new observation. This new observation data is from ERA5. Column 3: RMSE after including a pseudo new observation. Column 4: elevation map of the considered regions. The results are based on a year-long (Year 2021) evaluation.

378 For the case of Tibetan Plateau (first row), where the terrain is highly complex,  
 379 with average elevations exceeding 4500 meters, we obtain a relatively high RMSE given  
 380 existing observation constrains, particularly in the central and eastern parts of the re-  
 381 gion. Adding a station in the high-error area significantly reduces the RMSE for a broad  
 382 range of the considered region, this impact is more pronounced here as compared to the  
 383 other two cases, highlighting the importance of observational constraints in areas with  
 384 complex terrain. Removing a station results in a noticeable increase in RMSE in the sur-  
 385 rounding areas. Similarly, the effect of station removal is more evident compared to the

386 other two cases, suggesting that the model heavily relies on the limited observational data  
 387 to constrain its estimates in this complex terrain. The loss of a station in a critical lo-  
 388 cation can greatly impact the model’s ability to capture the local temperature patterns.

389 For the case of Peal River Delta (second row), the terrain is characterized by a mix  
 390 of lowlands and hilly regions, with elevations ranging from 0 to 1000 meters. The orig-  
 391 inal RMSE is low overall, with some higher values in the central and northwest moun-  
 392 tain regions. Adding a station in the high-error area effectively reduces the RMSE. Mean-  
 393 while, removing a station leads to a hardly noticeable increase in RMSE in the surround-  
 394 ing areas.

395 For the case of North China Plain (third row), the northern part is featured by moun-  
 396 tainous terrains exceeding 1000 meters, and the southern part has flat topography and  
 397 homogeneous terrain. Adding a station in the central of southern plain area reduces the  
 398 RMSE significantly, as existing observations are either from the northern mountain ar-  
 399 eas, or is too far away. Same as previous case, removing a station has minimal impact  
 400 on the RMSE distribution.

401 To sum up, we discuss the application of the CLIN methodology to evaluate the  
 402 impact of observational data on state estimation uncertainty across three diverse regions.  
 403 It emphasizes the importance of strategic addition and removal of observational stations  
 404 in improving estimation accuracy, particularly in areas with complex terrain. The find-  
 405 ings highlight how existing observation constraints influence RMSE distribution, with  
 406 significant reductions observed when stations are added in high-error areas. Conversely,  
 407 removal of stations leads to increased RMSE, underscoring the model’s reliance on lim-  
 408 ited observational data. Overall, we provide valuable insights for optimizing the design  
 409 of observation networks, leading to a reduction in uncertainties and biases in weather  
 410 and climate analysis.

## 411 4 Conclusion

412 Accurate state estimation of Earth atmosphere marks a daunting task due to its  
 413 high-dimensionality and chaotic nature. We demonstrated the potential of deep gener-  
 414 ative models, specifically probabilistic diffusion models, in learning the inherent low-dimensional  
 415 statistical structure of atmospheric circulation from climate reanalysis and simulation  
 416 data. By leveraging this learned climatological prior, we developed a methodology named  
 417 CLIN (Climate Inpainting) to effectively infer weather states from partial observations.

418 For the case study of estimating 2 m temperature spatial patterns, the learned cli-  
 419 matological prior accurately reproduced the essential characteristics and patterns of the  
 420 training data at both grid-scale and field-scale. This learned prior effectively captured  
 421 multi-scale climate patterns, providing regularization and stability to the state estima-  
 422 tion task.

423 Combining the learned climatological prior with station observations, CLIN yielded  
 424 strong posterior estimates of 2 m temperature spatial patterns. The estimates were spa-  
 425 tially coherent, well-constrained by observations, and provided reliable uncertainty quan-  
 426 tification. Regions near observation stations exhibited low ensemble variance, indicat-  
 427 ing high confidence in the estimates, while distant regions showed relatively higher en-  
 428 semble variance. The high spread-skill correlation confirmed that the ensemble spread  
 429 was a good predictor of the model’s estimation skill.

430 Moreover, CLIN allowed us to quantify the value of each observation station in re-  
 431 ducing state estimation uncertainty. By adding or removing stations and evaluating the  
 432 impact on the estimation error, we demonstrated the potential of this approach in guid-  
 433 ing the design of optimal observation networks.

Our study showcases the power of deep generative models in extracting and utilizing the information produced by the chaotic evolution of the climate system. The proposed CLIN methodology opens up new opportunities for data-driven weather state estimation, potentially complementing traditional data assimilation approaches.

Future work could focus on extending CLIN to handle indirect observations (i.e., remote sensing) and multiple interdependent variables, incorporating temporal dynamics, and adapting to long-term climate trends. Addressing the computational demands and data requirements of diffusion models is another important direction for making this approach more practical and accessible.

In conclusion, this study demonstrates the immense potential of deep generative models in advancing climate data exploration and tackling complex inference tasks in atmospheric sciences. By learning the intrinsic statistical structure of the climate system, these models can effectively bridge the gap between sparse observations and complete weather state estimates, paving the way for more accurate and efficient climate monitoring and prediction.

## 5 Data Availability

The ERA5 reanalysis data are obtained from the Copernicus Climate Change Service (C3S) Climate Data Store (CDS), accessible at <https://cds.climate.copernicus.eu/>.

The FGOALS model data are obtained from the Coupled Model Intercomparison Project Phase 6 (CMIP6), hosted by the Program for Climate Model Diagnosis and Intercomparison (PCMDI) at Lawrence Livermore National Laboratory (LLNL), accessible at <https://pcmdi.llnl.gov/CMIP6/>.

The observational data are freely available for download from the following website: <http://www.ncdc.noaa.gov/oa/ncdc.html>. The site information used in this study was obtained from the China Meteorological Data Network, hosted by the China National Meteorological Science Data Center (NMDC), accessible at <http://data.cma.cn/>.

## 6 Open Research

Model configuration, analysis scripts, data files used for this study will be publicly available upon accept of the work.

## Acknowledgments

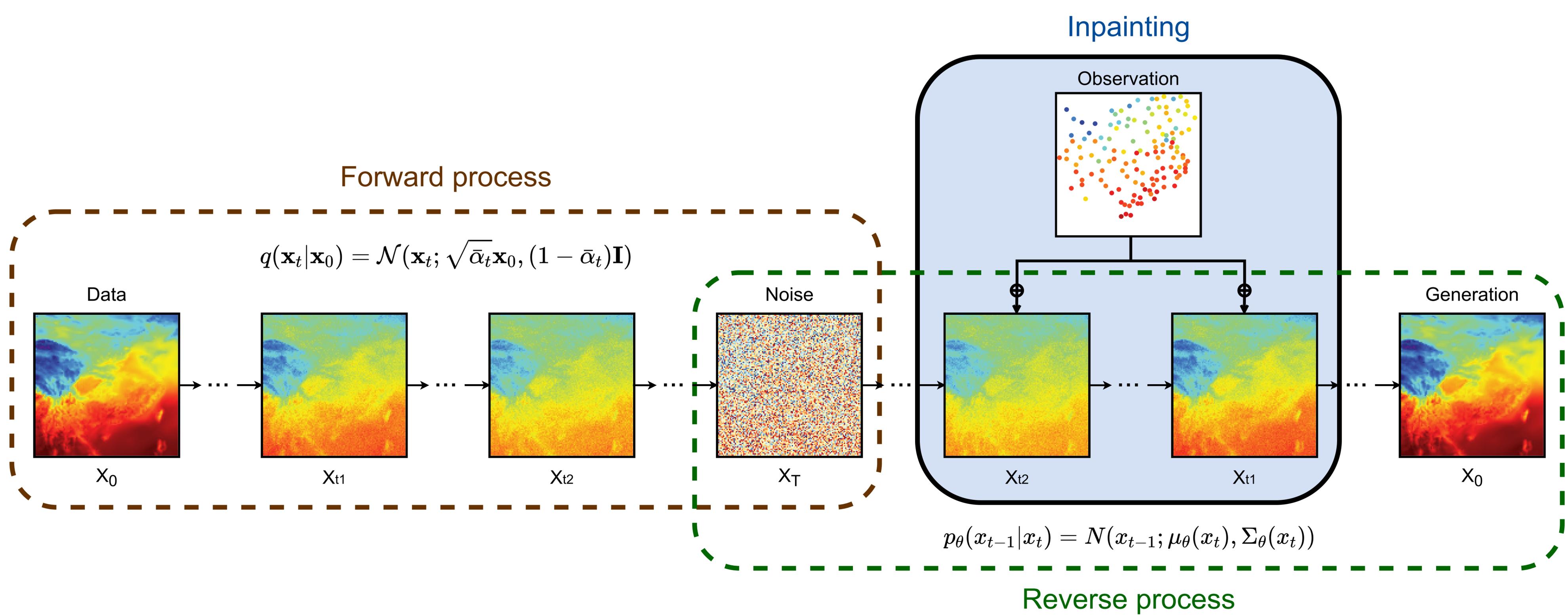
This research is supported by National Key R&D Program of China (Grant NoS. 2023YFC3007700 and 2023YFC3007705). We appreciate the insightful discussions with Dr. Niklas Boers from Technical University of Munich, and with Dr. Bin Wang and Dr. Juanjuan Liu from Chinese Academy of Science.

## References

- Bao, Q., Liu, Y., Wu, G., He, B., Li, J., Wang, L., ... others (2020). Cas fgoals-f3-h and cas fgoals-f3-l outputs for the high-resolution model intercomparison project simulation of cmip6. *Atmospheric and Oceanic Science Letters*, 13(6), 576–581.
- Carrassi, A., Bocquet, M., Bertino, L., & Evensen, G. (2018). Data assimilation in the geosciences: An overview of methods, issues, and perspectives. *Wiley Interdisciplinary Reviews: Climate Change*, 9(5), e535.
- Dhariwal, P., & Nichol, A. (2021). Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34, 8780–8794.

- 478 Eyring, V., Bony, S., Meehl, G. A., Senior, C. A., Stevens, B., Stouffer, R. J., &  
 479 Taylor, K. E. (2016). Overview of the coupled model intercomparison project  
 480 phase 6 (cmip6) experimental design and organization. *Geoscientific Model*  
 481 *Development*, 9(5), 1937–1958.
- 482 Ghil, M. (2020). Hilbert problems for the climate sciences in the 21st century–20  
 483 years later. *Nonlinear Processes in Geophysics*, 27(3), 429–451.
- 484 Gilpin, W. (2024). Generative learning for nonlinear dynamics. *Nature Reviews*  
 485 *Physics*, 1–13.
- 486 Hersbach, H., Bell, B., Berrisford, P., Hirahara, S., Horányi, A., Muñoz-Sabater, J.,  
 487 ... others (2020). The era5 global reanalysis. *Quarterly Journal of the Royal*  
 488 *Meteorological Society*, 146(730), 1999–2049.
- 489 Ho, J., Jain, A., & Abbeel, P. (2020). Denoising diffusion probabilistic models. *Ad-*  
 490 *vances in neural information processing systems*, 33, 6840–6851.
- 491 Holton, J. R., & Hakim, G. J. (2012). *An introduction to dynamic meteorology*. Aca-  
 492 *ademic press*.
- 493 Kadow, C., Hall, D. M., & Ulbrich, U. (2020). Artificial intelligence reconstructs  
 494 missing climate information. *Nature Geoscience*, 13(6), 408–413.
- 495 Kanngießer, F., & Fiedler, S. (2024). “seeing” beneath the clouds—machine-  
 496 learning-based reconstruction of north african dust plumes. *AGU Advances*,  
 497 5(1), e2023AV001042.
- 498 Kingma, D., Salimans, T., Poole, B., & Ho, J. (2021). Variational diffusion models.  
 499 *Advances in neural information processing systems*, 34, 21696–21707.
- 500 Law, K., Stuart, A., & Zygalakis, K. (2015). Data assimilation. *Cham, Switzerland:*  
 501 *Springer*, 214, 52.
- 502 Lugmayr, A., Danelljan, M., Romero, A., Yu, F., Timofte, R., & Van Gool, L.  
 503 (2022). Repaint: Inpainting using denoising diffusion probabilistic models.  
 504 In *Proceedings of the ieee/cvf conference on computer vision and pattern recog-*  
 505 *nition* (pp. 11461–11471).
- 506 Nai, C., Pan, B., Chen, X., Tang, Q., Ni, G., Duan, Q., ... Liu, X. (2024). Reli-  
 507 able precipitation nowcasting using probabilistic diffusion models. *Environmental*  
 508 *Research Letters*.
- 509 Palmer, T. (2017). The primacy of doubt: Evolution of numerical weather prediction  
 510 from determinism to probability. *Journal of Advances in Modeling Earth Sys-*  
 511 *tems*, 9(2), 730–734.
- 512 Pan, B., Anderson, G. J., Goncalves, A., Lucas, D. D., Bonfils, C. J., Lee, J., ...  
 513 Ma, H.-Y. (2021). Learning to correct climate projection biases. *Journal of*  
 514 *Advances in Modeling Earth Systems*, 13(10), e2021MS002509.
- 515 Pan, B., Wang, L.-Y., Zhang, F., Duan, Q., Li, X., Pan, X., ... others (2023). Prob-  
 516 abilistic diffusion model for stochastic parameterization—a case example of  
 517 numerical precipitation estimation. *Authorea Preprints*.
- 518 Pan, X., Guo, X., Li, X., Niu, X., Yang, X., Feng, M., ... others (2021). National  
 519 tibetan plateau data center: promoting earth system science on the third pole.  
 520 *Bulletin of the American Meteorological Society*, 102(11), E2062–E2078.
- 521 Schneider, T., Lan, S., Stuart, A., & Teixeira, J. (2017). Earth system modeling  
 522 2.0: A blueprint for models that learn from observations and targeted high-  
 523 resolution simulations. *Geophysical Research Letters*, 44(24), 12–396.
- 524 Sohl-Dickstein, J., Weiss, E., Maheswaranathan, N., & Ganguli, S. (2015). Deep  
 525 unsupervised learning using nonequilibrium thermodynamics. In *International*  
 526 *conference on machine learning* (pp. 2256–2265).
- 527 Song, Y., Garg, S., Shi, J., & Ermon, S. (2020). Sliced score matching: A scalable  
 528 approach to density and score estimation. In *Uncertainty in artificial intelli-*  
 529 *gence* (pp. 574–584).
- 530 Song, Y., Sohl-Dickstein, J., Kingma, D. P., Kumar, A., Ermon, S., & Poole, B.  
 531 (2020). Score-based generative modeling through stochastic differential equa-  
 532 tions. *arXiv preprint arXiv:2011.13456*.

- 533 Toth, Z., Talagrand, O., Candille, G., & Zhu, Y. (2003). Probability and ensemble  
534 forecasts. *Forecast verification: A practitioner's guide in atmospheric science*,  
535 *137*, 163.
- 536 Wang, B., Zou, X., & Zhu, J. (2000). Data assimilation and its applications. *Pro-*  
537 *ceedings of the National Academy of Sciences*, *97*(21), 11143–11144.
- 538 Zhang, G., Ji, J., Zhang, Y., Yu, M., Jaakkola, T. S., & Chang, S. (2023). Towards  
539 coherent image inpainting using denoising diffusion implicit models.



Mean

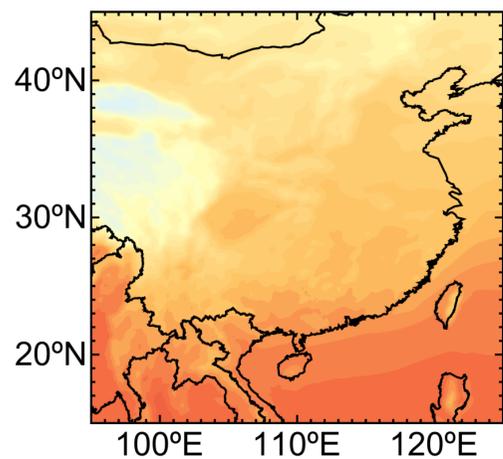
Variance

Skewness

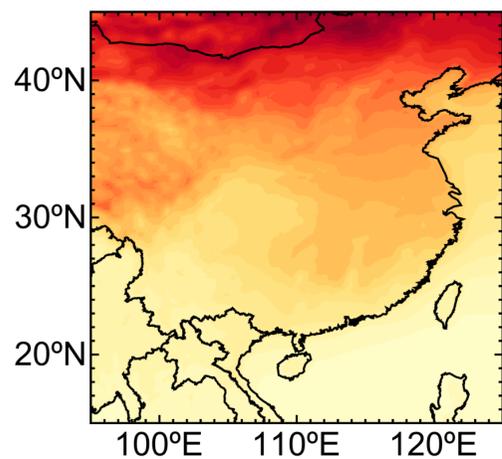
Min

Max

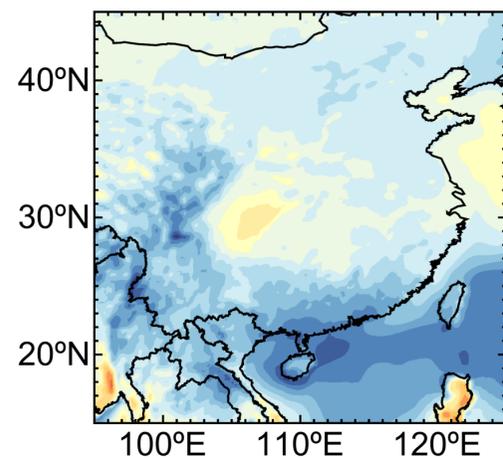
ERA5



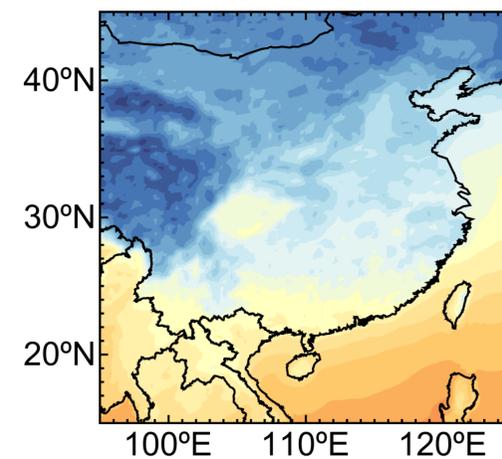
corr = 0.999 RMSE = 0.07°C



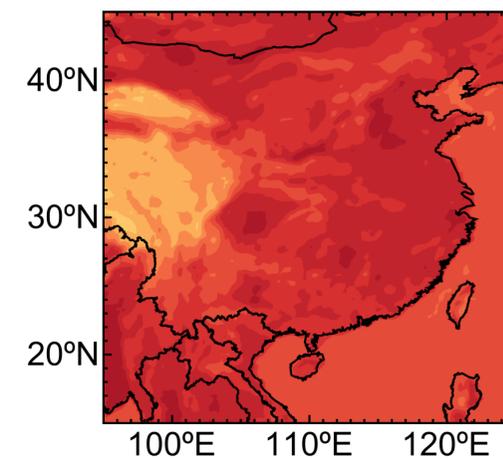
corr = 0.999 RMSE = 0.35(°C)²



corr = 0.999 RMSE = 0.004

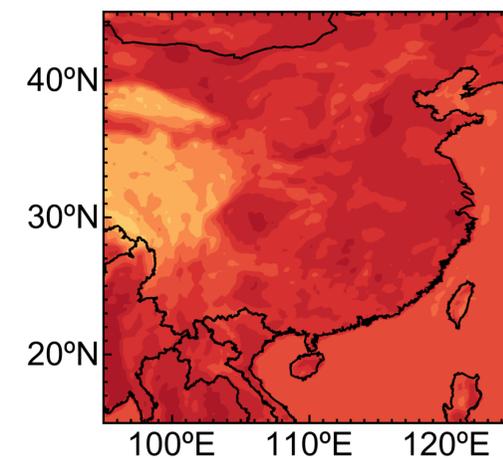
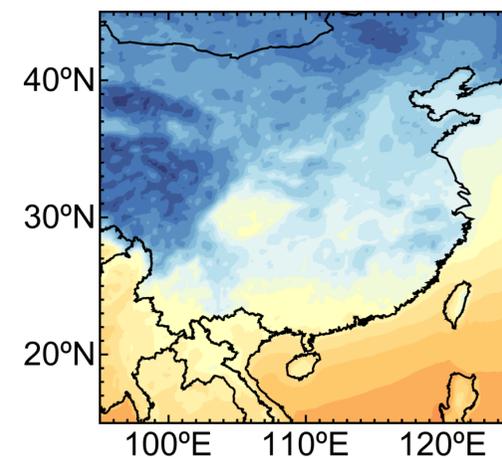
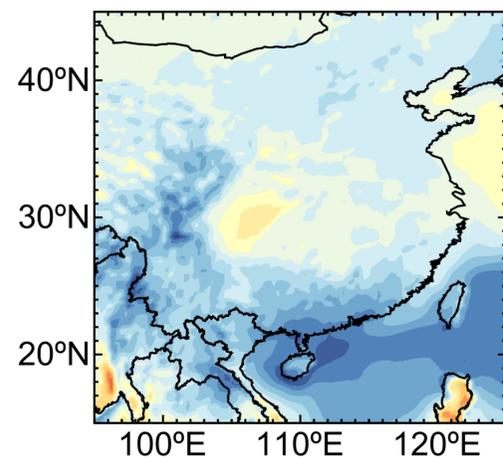
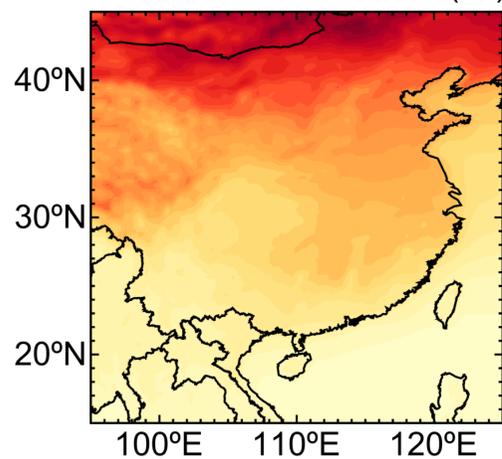
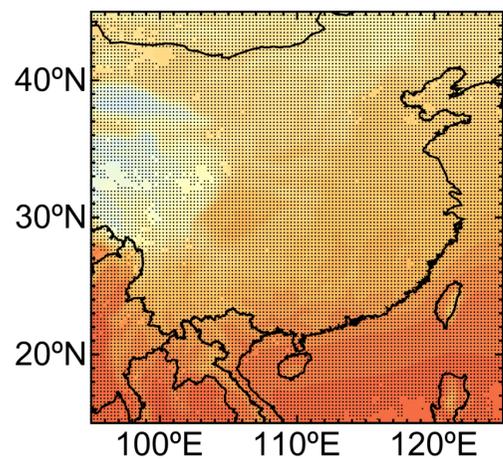


corr = 0.997 RMSE = 0.12°C

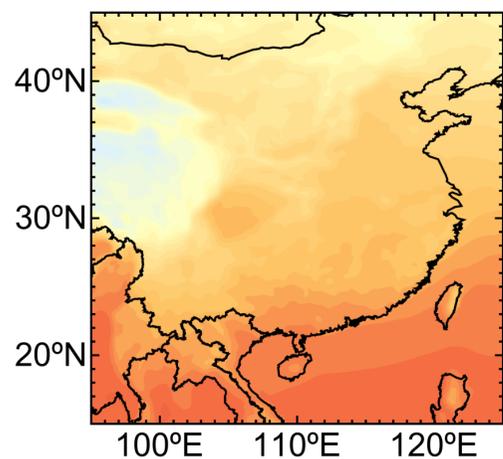


corr = 0.992 RMSE = 0.09°C

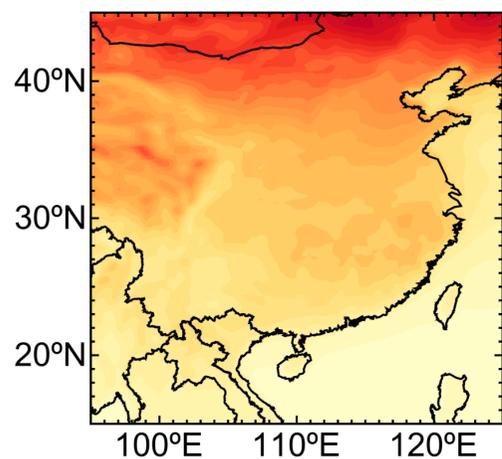
CLINERA5



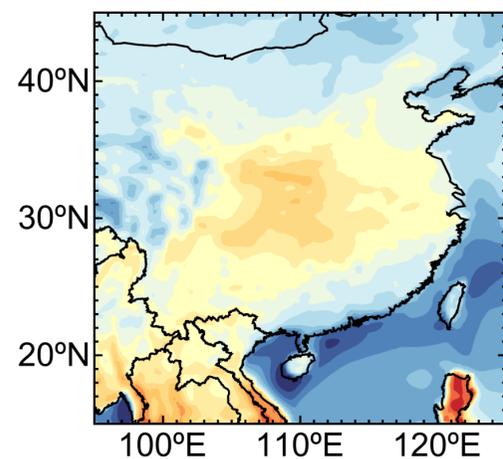
FGOALS



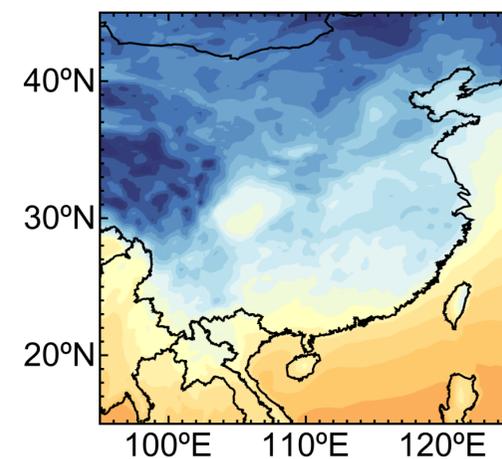
corr = 0.999 RMSE = 0.11°C



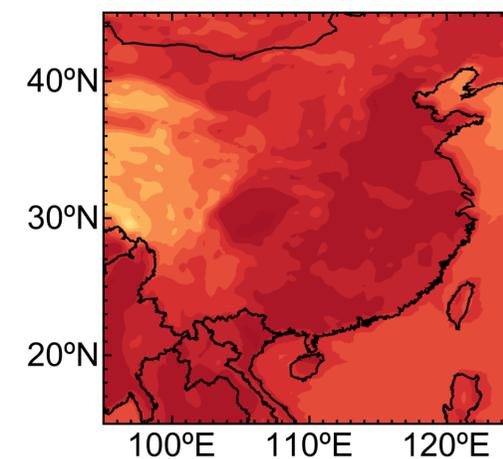
corr = 0.999 RMSE = 0.47(°C)²



corr = 0.997 RMSE = 0.033

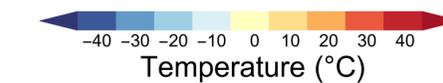
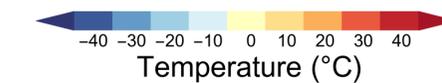
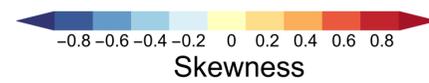
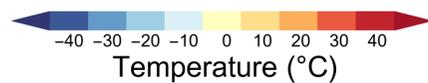
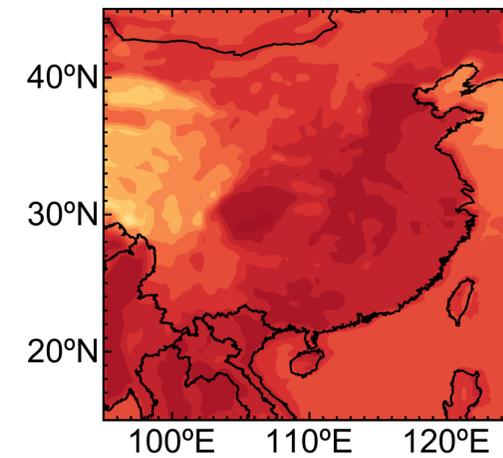
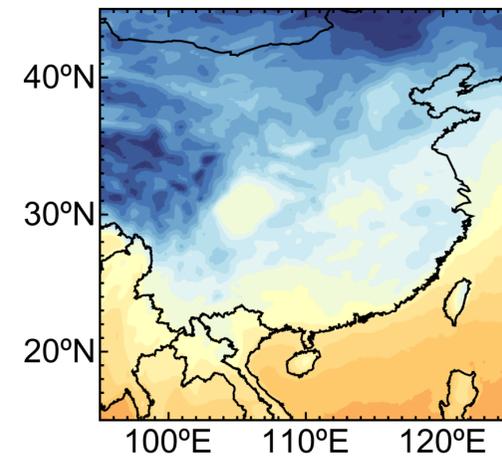
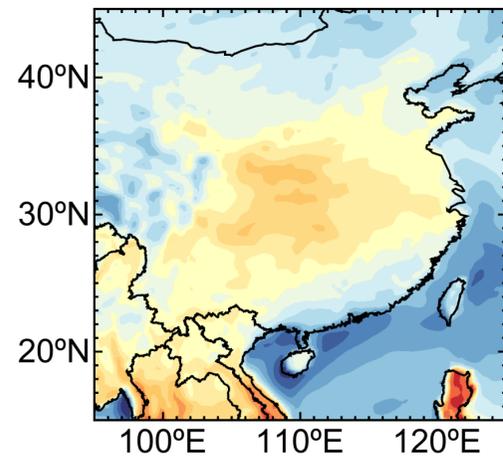
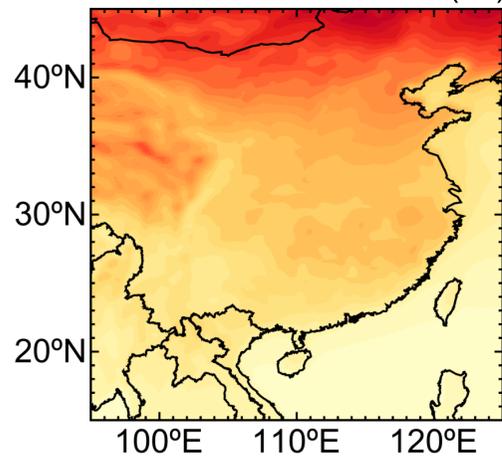
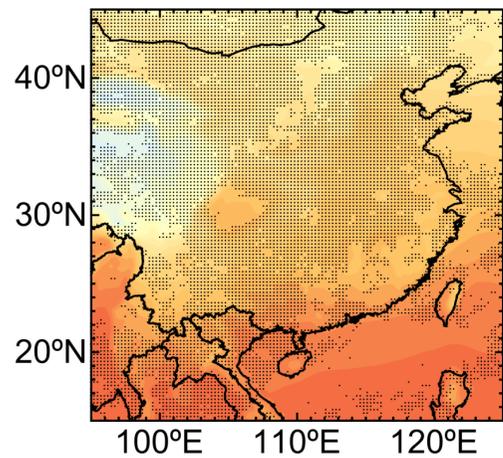


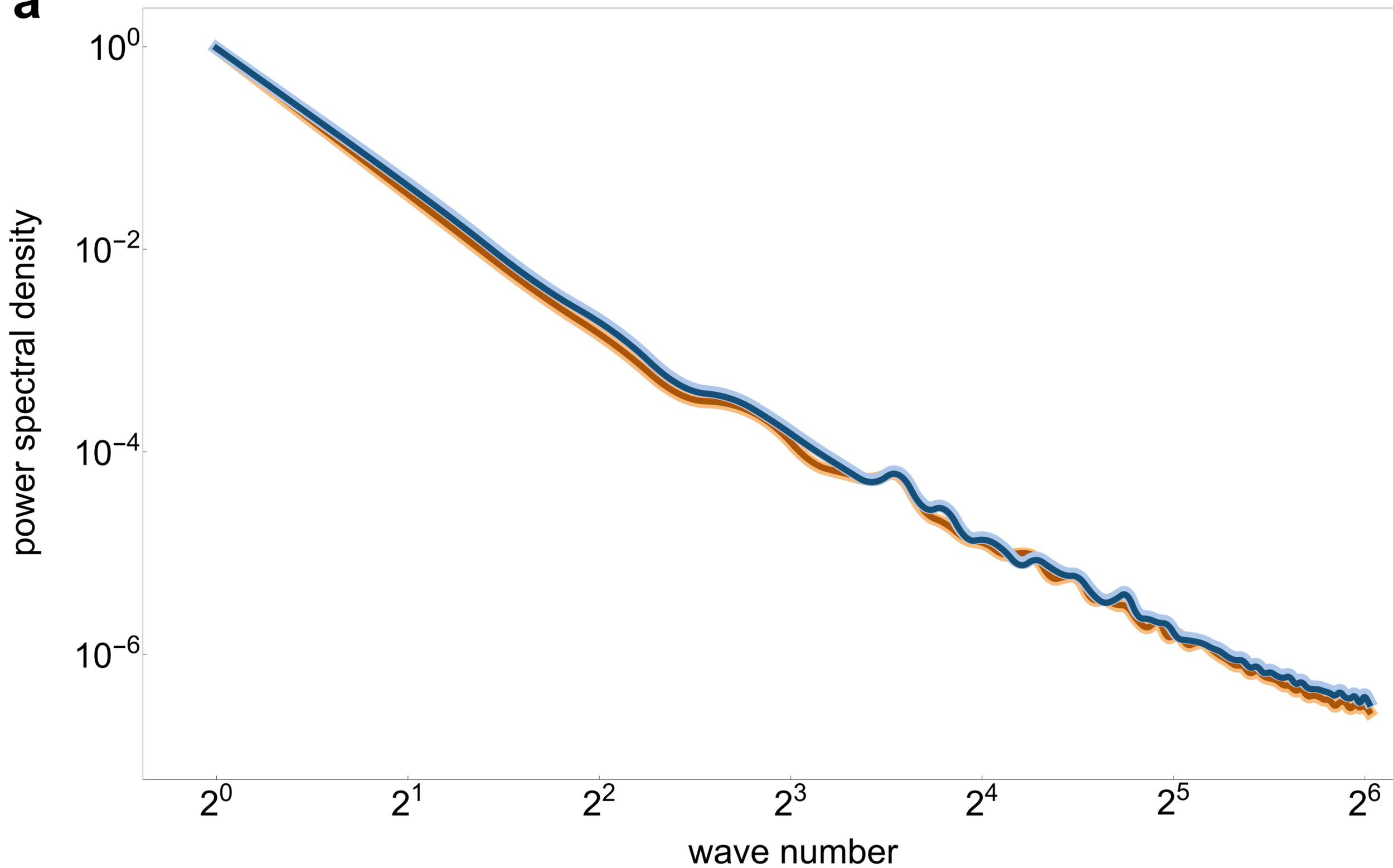
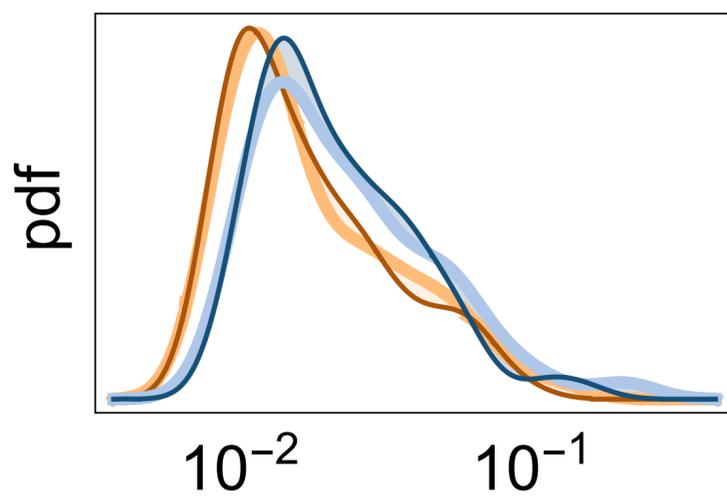
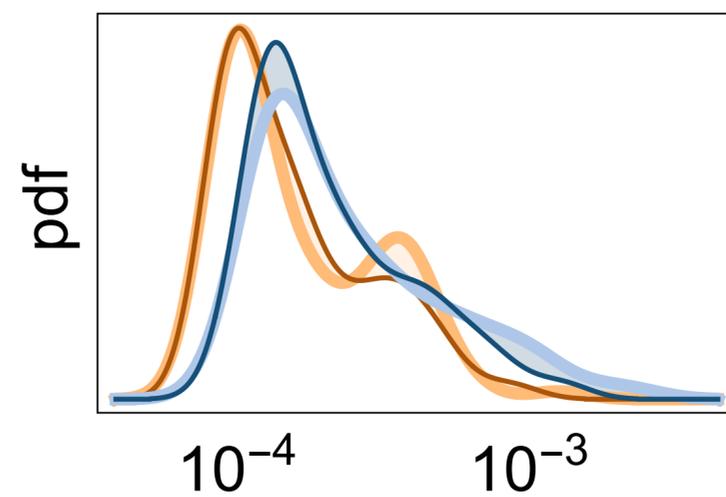
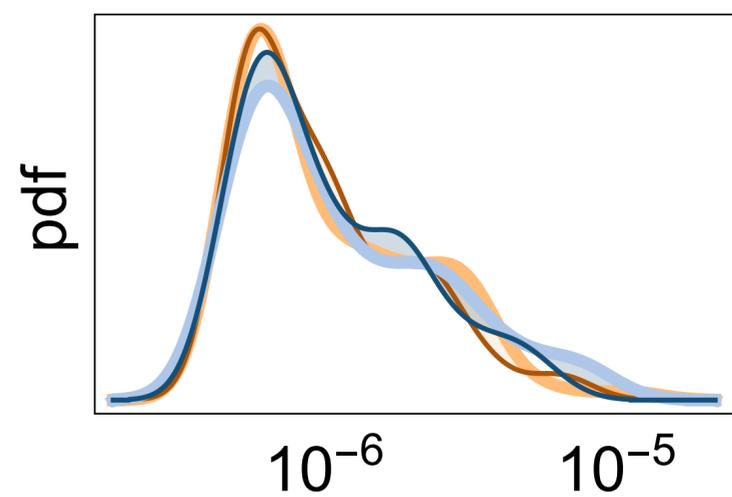
corr = 0.992 RMSE = 0.16°C



corr = 0.986 RMSE = 0.31°C

CLINFGOALS



**a****b**wave number  $2^1$ **c**wave number  $2^3$ **d**wave number  $2^5$ 

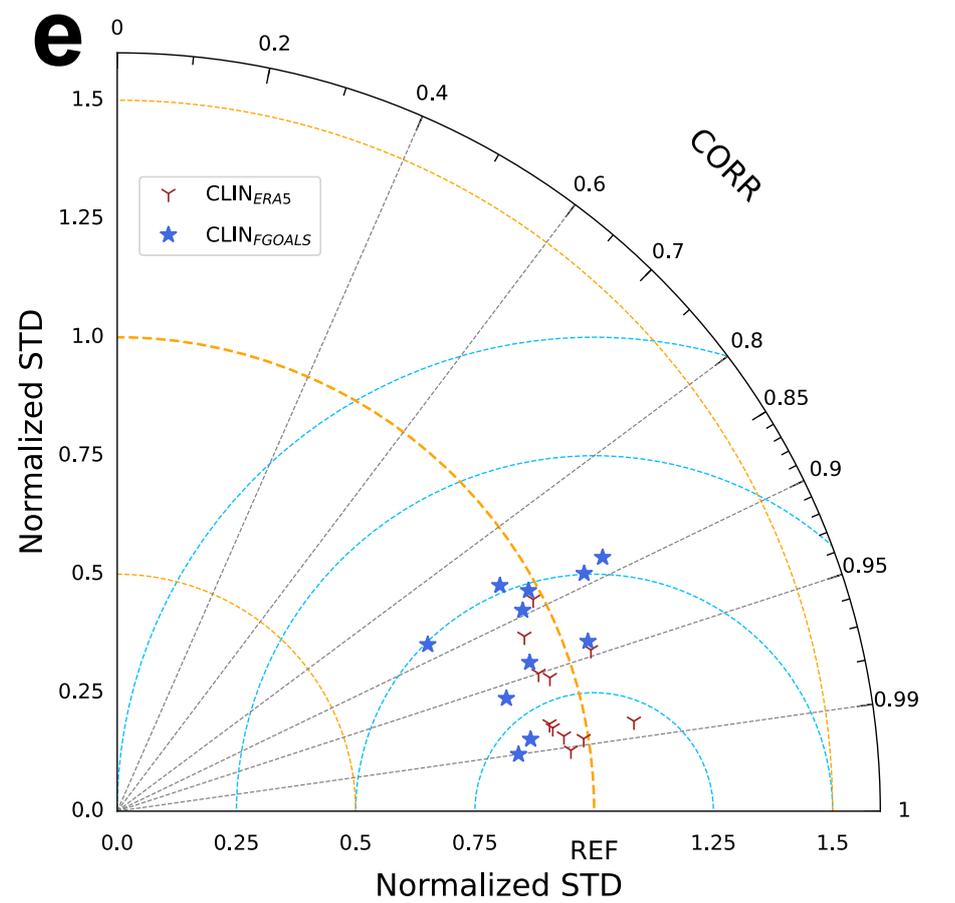
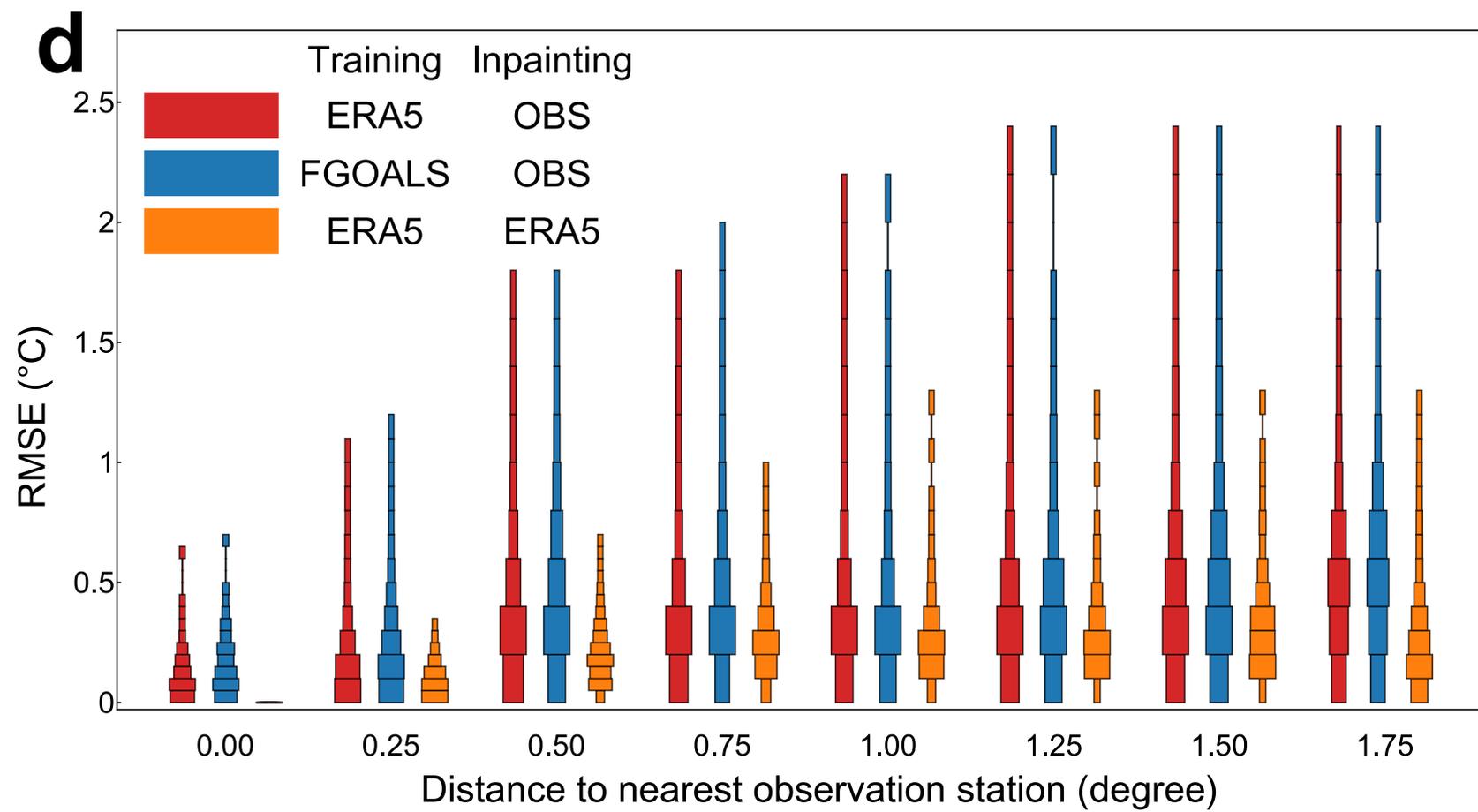
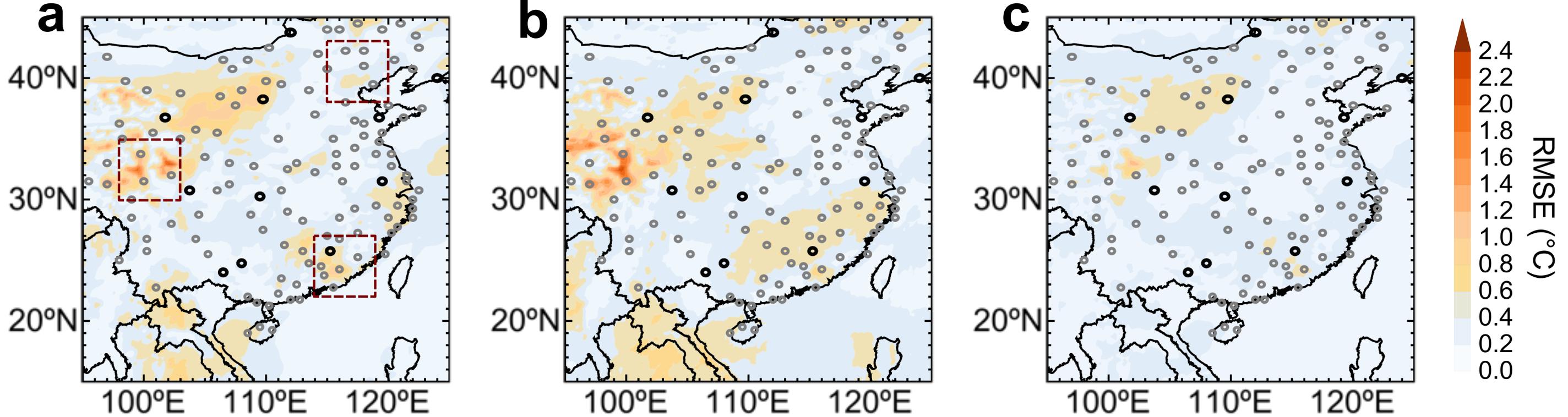
ERA5

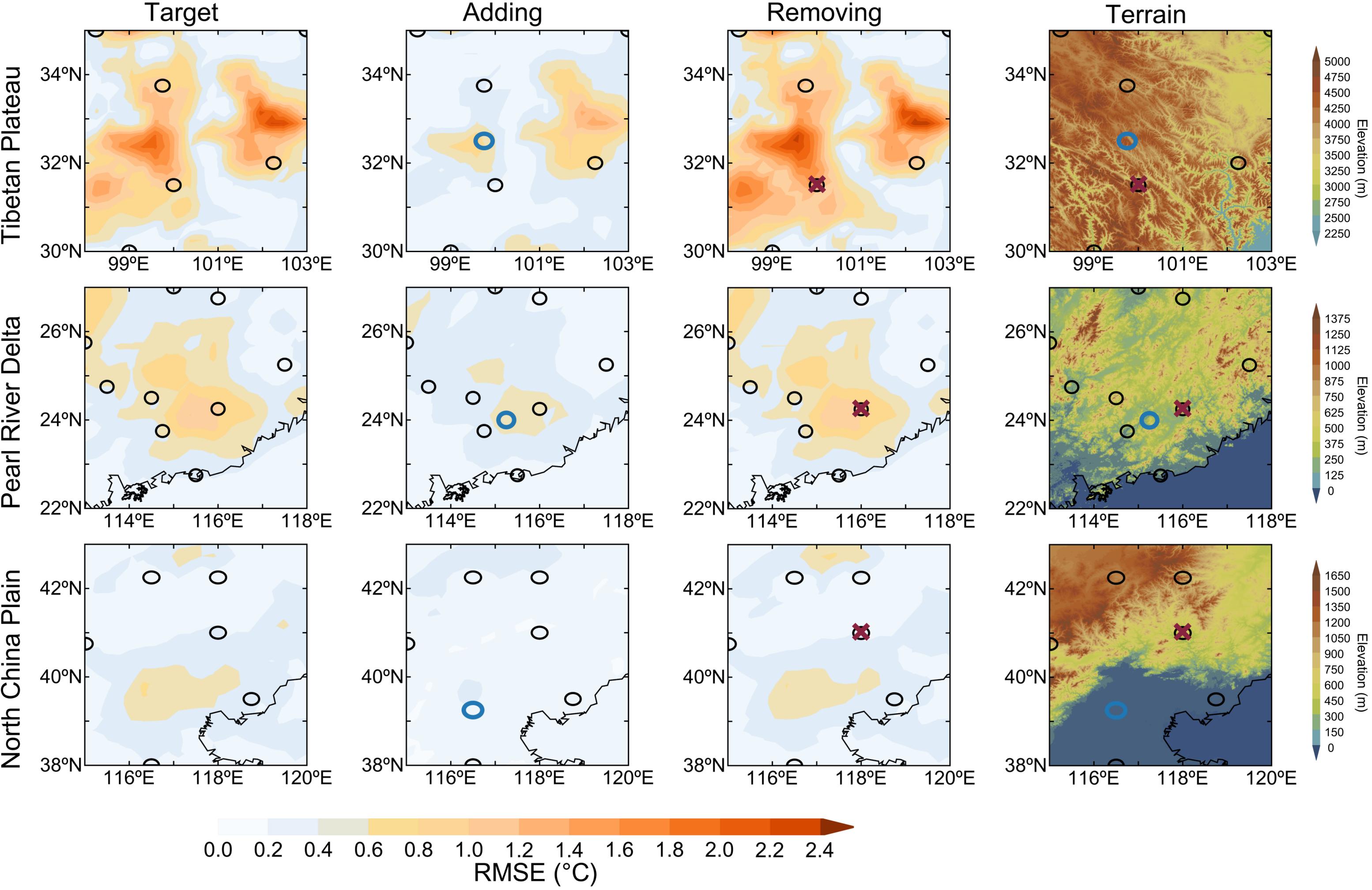
CLIN<sub>ERA5</sub>

FGOALS

CLIN<sub>FGOALS</sub>







# Learning to infer weather states using partial observations

Jie Chao<sup>1,2</sup>, Baoxiang Pan<sup>2</sup>, Quanliang Chen<sup>1</sup>, Shangshang Yang<sup>2,3</sup>, Jingnan Wang<sup>2,4</sup>, Congyi Nai<sup>2,5</sup>, Yue Zheng<sup>6</sup>, Xichen Li<sup>2</sup>, Huiling Yuan<sup>3</sup>, Xi Chen<sup>2</sup>, Bo Lu<sup>7</sup>, Ziniu Xiao<sup>2</sup>

<sup>1</sup>School of Atmospheric Sciences, Chengdu University of Information Technology, Sichuan, China

<sup>2</sup>Institute of Atmospheric Physics, Chinese Academy of Science, Beijing, China

<sup>3</sup>Key Laboratory of Mesoscale Severe Weather, Ministry of Education, and School of Atmospheric Sciences, Nanjing University, Jiangsu, China

<sup>4</sup>College of Computer, National University of Defense Technology, Hunan, China

<sup>5</sup>Institute of Geographic Sciences and Natural Resources Research, Chinese Academy of Sciences, Beijing, China

<sup>6</sup>Clustertech LTD, Hong Kong, China

<sup>7</sup>National Climate Center, China Meteorological Administration, Beijing, China

## Key Points:

- Deep generative model enables accurate spatial interpolation of weather variables from sparse observations.
- The model generates probabilistic weather estimates with reliable uncertainty quantification by combining learned priors and observations.
- The model quantifies the value of observations for reducing uncertainty, guiding optimal observation network design.

---

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

## Abstract

Accurate state estimation of the high-dimensional, chaotic Earth atmosphere marks a Sisyphean task, yet is indispensable for initiating weather forecast and gauging climate variability. While much effort is devoted to assimilating observations and forecasts to infer weather state, the inherent low-dimensional statistical structure in atmospheric circulation, shaped by geophysical laws and geographic boundaries, is underutilized as informative prior for state inference, or as reference for assessing representative of existing observations and planning new ones. We realize these potential by learning climatological distribution from climate reanalysis/simulation, using deep generative model. For a case study of estimating 2 m temperature spatial patterns, the learned distribution faithfully reproduces climatology statistics. A combination of the learned climatological prior with few station observations yields strong posterior of spatial pattern estimates, which are spatially coherent, faithful and adaptive to observation constraints, and uncertainty-aware. This allows us to evaluate each observation's value in reducing state estimation uncertainty, and guide optimal observation network design by pinpointing the most informative sites. Our study showcases how generative models can extract and utilize information produced in the chaotic evolution of climate system.

## Plain Language Summary

Accurate estimation of weather conditions across a large area is crucial but challenging due to the complex and chaotic nature of the atmosphere. Traditional methods rely on combining observations with forecasts, which can be computationally expensive and sensitive to model biases. We propose a new approach called Climate Inpainting (CLIN) that learns the inherent spatial patterns of the atmosphere from climate data using machine learning techniques. CLIN can effectively combine the learned patterns with limited observations to reconstruct complete spatial maps of weather variables, such as temperature. We demonstrate that CLIN can accurately reproduce the key spatial features and variability of temperature over East Asia. Moreover, CLIN can quantify the uncertainty in the estimated weather maps and evaluate the importance of each observation site in reducing the overall uncertainty. This information can guide the optimal design of weather station networks. Our approach showcases the potential of machine learning in utilizing the rich information contained in climate data to improve weather estimation and observation planning.

## 1 Introduction

The state of the Earth atmosphere, which concerns a broad range of socioeconomic sectors and the overall environment, is characterized by the spatial distribution of a specific set of physical properties, including temperature, pressure, wind speed and direction, density, concentration of water of different phases, composition of aerosol, greenhouse gas, etc (Holton & Hakim, 2012). To determine the atmosphere state at 50 km grid resolution requires estimating the value for all the above-mentioned physical properties at around  $\sim 10^7$  grids (Schneider et al., 2017). Doubling the resolution increases the total number of grids by a factor of 8. This high dimensionality poses a daunting challenge for monitoring the atmosphere (Ghil, 2020).

Current operational forecasting centers routinely update their atmosphere state estimates by combining multi-source observations and previous forecasts, so as to reboot weather forecast and gauge climate variability (Carrassi et al., 2018). Ground based observations offer direct meteorological measurements, yet come with limited spatial coverage and high maintenance cost. Remote sensing offers broader spatial coverage, yet is indirect and error prone, requiring careful calibration based on ground-based observations.

71 Deficiencies in observation render it an ill-posed task to estimate the state of the  
72 high-dimensional Earth atmosphere, calling for strong prior to achieve feasible solution.  
73 Forecasts from previous time steps are frequently applied to serve this mission, carry-  
74 ing information from previous step observations to the current step via a process-based  
75 model (Wang et al., 2000). As a result, the state estimation accuracy depends on an in-  
76 tricate interplay among model biases, background uncertainty, and observation error, which  
77 cannot be effectively disentangled or controlled (Law et al., 2015). Moreover, to provide  
78 multi-scale background information using forecasting models requires operational run  
79 of large ensemble high-resolution numerical simulations, which is prohibitively expen-  
80 sive and burdensome (Toth et al., 2003; Palmer, 2017).

81 Is there extra information source for inferring the state of the high-dimensional,  
82 chaotic Earth atmosphere? It turns out that, the inherent low-dimensional statistical struc-  
83 ture in atmospheric circulation, shaped by the underlying geophysical laws and quasi-  
84 static geographic boundaries, can serve as an informative prior for state inference. The  
85 Earth climate system, like any other chaotic system, is an information producer: it grad-  
86 ually reveals the characteristic structure of its phase space at ever-finer scales (Gilpin,  
87 2024). By identifying and parameterizing this characteristic structure, we can potentially  
88 bypass the curse of high dimensionality, and make more efficient use of limited obser-  
89 vations for the state inference task.

90 Some pioneering works have explored this direction, leveraging the inherent struc-  
91 ture of climate data to fill in missing observations and rebuild historical climate records.  
92 For instance, Kadow et al. (2020) developed a partial convolution method to reconstruct  
93 historical global temperature patterns based on partial observations and climate simu-  
94 lation. Kanngießer and Fiedler (2024) applied a similar methodology to restore the spa-  
95 tial extent of dust plumes in cloud-masked satellite images. Most of these practices con-  
96 sider deterministic models, which are designed for specific “reconstruction” problem con-  
97 figurations, yielding deterministic results regardless of whether observations can adequately  
98 constrain the estimation uncertainty. As a result, these methodologies generalize poorly  
99 to state inference tasks where the number or layout of observations change, fail to re-  
100 produce extremes or apply for scenarios where only limited observations are available.

101 A solution to these dilemmas is to shift from deterministic model to probabilistic  
102 model (B. Pan et al., 2021). Specifically, we prefer to build a probabilistic model that  
103 explicitly represents the inherent statistical structure of the atmosphere as revealed by  
104 climate observations or simulations. Thereafter, we hope to effectively and efficiently com-  
105 bine the learned climatological prior with incomplete observations, so as to obtain strong  
106 posterior of spatial pattern estimates. This problem setup poses two stringent require-  
107 ments on the underlying probabilistic model. First, the model must faithfully approx-  
108 imate the high-dimensional climatological distribution as generated by the chaotic evo-  
109 lution of climate dynamics. Second, the model must enable flexible probabilistic infer-  
110 ence, allowing us to efficiently obtain posterior atmospheric state estimates given arbi-  
111 trary observational constraints.

112 To fulfill these requirements, we resort to generative machine learning, in partic-  
113 ular, probabilistic diffusion models (Sohl-Dickstein et al., 2015; Ho et al., 2020; Song, Sohl-  
114 Dickstein, et al., 2020; Kingma et al., 2021). Probabilistic diffusion models learn to ap-  
115 proximate complex, high-dimensional probability distributions in an iterative manner,  
116 achieving unprecedented fitting capacity and controlling flexibility (B. Pan et al., 2023;  
117 Nai et al., 2024). To demonstrate the idea, we consider a case example of inferring the  
118 spatial pattern of 2 m temperature based on sparse observations from operational me-  
119 teorology stations. We learn probabilistic diffusion models to approximate the climato-  
120 logical distribution of 2 m temperature spatial patterns from climate reanalysis or simu-  
121 lation data. After carefully assessing the model’s ability to reproduce climatology, we  
122 develop tools to “inpaint” arbitrary observation constraints into the sample generation  
123 process, yielding probabilistic 2 m temperature spatial pattern estimates. Finally, we ap-

124 ply this methodology to evaluate each observation’s value in reducing state estimation  
 125 uncertainty, and guide optimal observation network design by pinpointing the most in-  
 126 formative sites.

## 127 2 Methodology

### 128 2.1 Data and problem setup

129 We consider the task of inferring the spatial pattern of 2 m temperature over East  
 130 Asia ( $15^\circ\text{N} - 45^\circ\text{N}, 95^\circ\text{E} - 125^\circ\text{E}$ ), using station observations covering  $\sim 1\%$  grids of  
 131 the considered region. To achieve this, we learn climatological distribution of 2 m tem-  
 132 perature spatial pattern using climate reanalysis or simulation data. The reanalysis data  
 133 are hourly,  $0.25^\circ$  2 m temperature data from the fifth-generation global climate and weather  
 134 reanalysis (ERA5) developed at European Centre for Medium-Range Weather Forecasts  
 135 (Hersbach et al., 2020, ECMWF). The simulation data are 3-hourly,  $0.25^\circ$  2 m temper-  
 136 ature historical simulation from the Flexible Global Ocean-Atmosphere-Land System Model  
 137 version f3-H (Bao et al., 2020, FGOALS-f3-H), which participates in the sixth phase of  
 138 the Coupled Model Intercomparison Project (Eyring et al., 2016, CMIP6). The station  
 139 observation data are obtained from the Chinese National Climatic Data Center (X. Pan  
 140 et al., 2021).

Formally, we denote the spatial pattern of 2 m temperature for the target region  
 as  $\mathbf{x}$ , which is a  $120 \times 120$  dimensional random variable here. Our objective is to ap-  
 proximate the distribution of  $\mathbf{x}$ , based on large number of samples from climate reanal-  
 ysis or simulation:

$$p_{\theta^*} = \arg \max_{p_\theta} \sum \log p_\theta(\mathbf{x}) \quad (1)$$

141 Here  $p_\theta$  is parameterized probability density function approximator,  $\theta^*$  is the optimal  
 142 parameter, optimized by maximizing the overall likelihood of  $p_\theta$  assigned to the train-  
 143 ing samples.

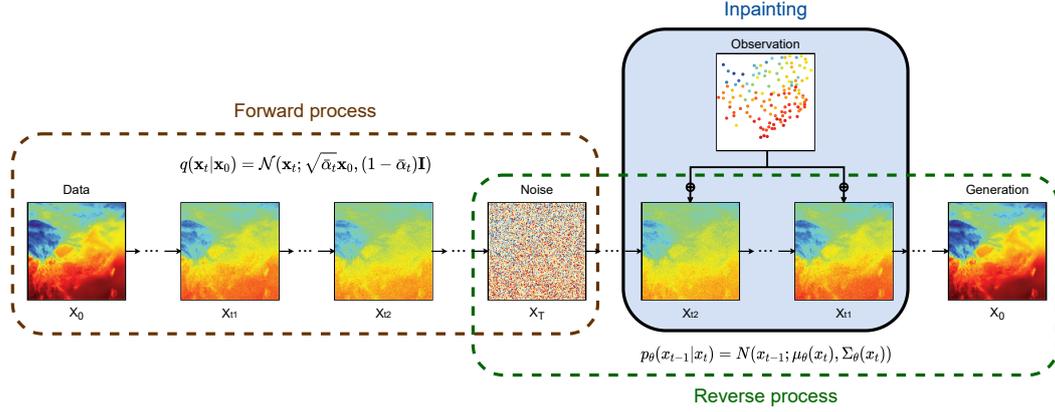
144 Given  $p_{\theta^*}$  and sparse observations, we need to provide probabilistic estimates of  
 145 2 m temperature spatial patterns, i.e.,  $p_{\theta^*}(\mathbf{x}|\mathbf{x} \odot \mathbf{m})$ . Here,  $\odot$  is dot product,  $\mathbf{m}$  is ob-  
 146 servation mask, with value 1/0 denoting the existence/absence of observations for each  
 147 geogrid.  $p_{\theta^*}(\mathbf{x}|\mathbf{x} \odot \mathbf{m})$  should yield samples that are spatially coherent and faithful to  
 148 observational constraints. Also,  $p_{\theta^*}(\mathbf{x}|\mathbf{x} \odot \mathbf{m})$  should offer accurate uncertainty quan-  
 149 tification. For instance, geogrids close to observation stations should typically have low  
 150 state estimate uncertainties, while distant ones have high uncertainties. Finally, we pre-  
 151 fer  $p_{\theta^*}(\mathbf{x}|\mathbf{x} \odot \mathbf{m})$  to be adaptive to changes in observation configurations, such as the  
 152 abortion or inclusion of observation stations, or rearrangement of station network lay-  
 153 out. Below we illustrate how to achieve these requirements using the proposed method-  
 154 ology.

### 155 2.2 Learning climatology with probabilistic diffusion model

156 We elucidate how to learn climatological distribution of the target random vari-  
 157 able using probabilistic diffusion model, thereafter leverage this learned prior for the in-  
 158 ference task (Sec. 2.3). For clarity, we only cover key steps necessary for establishing our  
 159 methodology. Details can be found in the literature referenced through the description.

160 To approximate a target distribution using probabilistic diffusion model, we train  
 161 a series of deep neural networks that can be chained to establish bijective mapping be-  
 162 tween the target distribution and a prior distribution (Sohl-Dickstein et al., 2015; Ho et  
 163 al., 2020). Specifically, we define the following Gaussian process:

$$q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{(1 - \beta_t)}\mathbf{x}_{t-1}, \beta_t\mathbf{I}) \quad (2)$$



**Figure 1.** Overview of the Climate Inpainting (CLIN) methodology. A pre-defined forward Gaussian process (left) turns distribution of target climate variable into a prior distribution, i.e., standard Gaussian. A learned reverse Gaussian process (right) turns the prior distribution into the distribution of the target climate variable. We “inpaint” sparse observations throughout the reverse Gaussian process (right top), so as to obtain spatial pattern estimates of the target variable.

164 Here  $p(\mathbf{x}_0) = p(\mathbf{x})$ , which is the target distribution;  $p(\mathbf{x}_T)$  is the prior distribution; we  
 165 bridge  $\mathbf{x}_0$  and  $\mathbf{x}_T$  using  $\mathbf{x}_{t \in [1, T]}$ , which are latent variables with increasing noise level;  
 166  $\mathcal{N}$  is Gaussian distribution;  $\mathbf{I}$  is identity matrix;  $\beta_t$  is diffusion coefficient, which is pre-  
 167 defined so that, give large enough  $T$ ,  $p(\mathbf{x}_T | \mathbf{x}_0)$  is drawn close to  $p(\mathbf{x}_T)$ , which is  $\mathbf{x}_0$  ag-  
 168 nostic. This setup offers analytical solution for  $p(\mathbf{x}_{t+\tau} | \mathbf{x}_t), \forall \tau \in [0, T - t], t \in [0, T]$ ,  
 169 facilitating convenient inference as detailed in Sec. 2.3.

To achieve generative modeling, we reverse Eq. 2 using the following variation distributions:

$$p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \mu_\theta, \Sigma_\theta) \quad (3)$$

170 Here  $\Sigma_\theta$  is represented as an interpolation between its analytical lower and upper bound  
 171 (Dhariwal & Nichol, 2021);  $\mu_\theta$  can be optimized by maximizing the variational lower bound  
 172 (ELBO) on the log-likelihood of the training samples (Sohl-Dickstein et al., 2015; Kingma  
 173 et al., 2021). In practice, we represent  $\mu_\theta$  as function of neural network parameteriza-  
 174 tion for  $\nabla p(\mathbf{x}_t | \mathbf{x}_0)$ , which is known as the *score function* (Song, Garg, et al., 2020; Song,  
 175 Sohl-Dickstein, et al., 2020). This simplifies the ELBO objective function to the follow-  
 176 ing form:

$$L = \mathbb{E}_{t \in [1, T], \mathbf{x}_0 \sim p(\mathbf{x}_0)} \|\nabla p(\mathbf{x}_t | \mathbf{x}_0) - \epsilon_\theta\|^2 \quad (4)$$

177 Here  $\epsilon_\theta$  is a neural network parameterization for  $\nabla p(\mathbf{x}_t | \mathbf{x}_0)$ . Given the trained score es-  
 178 timates, we can derive  $p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \mu_\theta, \Sigma_\theta)$  and sample it, starting with  $p(\mathbf{x}_T)$ ,  
 179 ending with  $p(\mathbf{x}_0)$ .

### 2.3 CLIN: inferring weather states using partial observations

180  
 181 We combine the learned climatology prior with station observations to infer the pos-  
 182 terior probability distribution of the target variable, using a *repainting* methodology (Lugmayr  
 183 et al., 2022; Zhang et al., 2023). Specifically, given a pre-trained diffusion model that se-  
 184 quentially applies  $p_{\theta^*}(\mathbf{x}_t | \mathbf{x}_{t+1}) = \mathcal{N}(\mathbf{x}_t; \mu_{\theta^*}, \Sigma_{\theta^*})$  to transform  $p(\mathbf{x}_T)$  to  $p(\mathbf{x}_0)$ , within  
 185 a pre-selected time window of  $\Omega$ , for grid points where we have observations, we replace

186 values of  $\mathbf{x}_t$  with observations noisified to time step  $t$ , by sampling  $p(\mathbf{x}_t \odot \mathbf{m} | \mathbf{x}_0 \odot \mathbf{m})$ .  
 187 This replacement does not consider the generated parts of  $\mathbf{x}_t$ , therefore, the observations  
 188 could not explicitly constrain the variability of unobserved parts.

189 To address this issue, for any  $t \in \Omega$ , after the replacement, instead of progress-  
 190 ing to  $t - 1$  directly, we rewind to time step  $t - \tau$  by sampling  $p(\mathbf{x}_{t-\tau} | \mathbf{x}_t)$ . We there-  
 191 after repeat the denoising steps from  $t - \tau$  to  $t$  for  $k$  rounds, and carry out observation  
 192 replacement for  $\mathbf{x}_t$  at each round. This allows us to jointly modify both observed and  
 193 unobserved regions throughout the denoising steps, yielding generated samples that are  
 194 spatially coherent, faithful and adaptive to observation constraints, and uncertainty-aware.  
 195 This methodology is referred to as *inpainting*, we hence name our methodology as CLIN,  
 196 short for Climate Inpainting. A formal algorithm description is given below. Details for  
 197 data processing, neural network architecture, hyperparameters for training and inference,  
 198 are given in Supporting Information.

---

**Algorithm 1** CLIN

---

**Require:** trained diffusion model  $p_{\theta^*}$ , observations  $\mathbf{x}_0 \odot \mathbf{m}$ , repainting time step set  $\Omega$ ,  
 rewinding step  $\tau$ , rewinding round  $K$   
**Ensure:** observation constrained, spatially coherent sample  $\mathbf{x}_0$

- 1: Initialize  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 2: **for**  $t = T - 1, \dots, 1$  **do**
- 3:      $\mathbf{x}_t \sim p_{\theta^*}(\mathbf{x}_t | \mathbf{x}_{t+1})$  ▷ Reverse sampling
- 4:     **if**  $t \in \Omega$  **then:**
- 5:         **for**  $k = 1, \dots, K$  **do**
- 6:              $\mathbf{x}_t^{\text{obser}} \sim p(\mathbf{x}_t \odot \mathbf{m} | \mathbf{x}_0 \odot \mathbf{m})$
- 7:              $\mathbf{x}_t \leftarrow \mathbf{x}_t \odot (\mathbf{I} - \mathbf{m}) + \mathbf{x}_t^{\text{obser}}$  ▷ Condition on observations
- 8:              $\mathbf{x}_{t+\tau} \sim p(\mathbf{x}_{t+\tau} | \mathbf{x}_t)$  ▷ Rewind in time by  $\tau$  steps
- 9:             **for**  $i = t + \tau - 1, \dots, t$  **do**
- 10:                  $\mathbf{x}_i \sim p_{\theta^*}(\mathbf{x}_i | \mathbf{x}_{i+1})$  ▷ Reverse sampling within a rewinding round
- 11:             **end for**
- 12:         **end for**
- 13:     **end if**
- 14: **end for**
- 15: **return**  $\mathbf{x}_0$

---

199 **3 Results**

200 The accuracy for state estimation depends on 1) how well we can approximate the  
 201 climatological distribution, and 2) based on a learned climatological prior, how well we  
 202 can combine it with limited observations to obtain probabilistic state estimates. Below  
 203 we assess model’s performance for these two aspects (Sec. 3.1 and 3.2). We further em-  
 204 ploy the model to quantify the extent to which observations reduce uncertainty in state  
 205 estimation, offering insights for optimal observation design (Sec. 3.3).

206 **3.1 Climatology**

207 We compare grid-scale and field-scale statistics of 10,000 reference/generated sam-  
 208 ples to evaluate how well the probabilistic diffusion models reproduce their training data’s  
 209 climatology. Two models trained with climate reanalysis (ERA5) and historical climate  
 210 simulation (FGOALS) data, hereafter referred to as  $\text{CLIN}_{\text{ERA5}}$  and  $\text{CLIN}_{\text{FGOALS}}$ , are  
 211 deployed and evaluated.

212 The grid-scale assessment considers the mean, variance, skewness, minimum, and  
 213 maximum of climatological distribution at each grid (Fig. 2). These statistics from ERA5  
 214 (Fig. 2 Row 1) and FGOALS (Fig. 2 Row 3) generally agree well, due to shared constraints  
 215 from geophysical laws and geographic boundaries. The key spatial patterns are the lat-  
 216 itudinal gradient, the influence of topography (e.g., the Tibetan Plateau), and the land-  
 217 sea contrast, which are most evident in the mean, minimum and maximum maps. The  
 218 variance and skewness maps reveal more regional variations. A notable discrepancy is  
 219 that, compared to ERA5, FGOALS tends to hold larger skewness for most of the land  
 220 regions in Southern China and Philippine Island, implying a more frequent present of  
 221 high 2 m temperature for these regions.

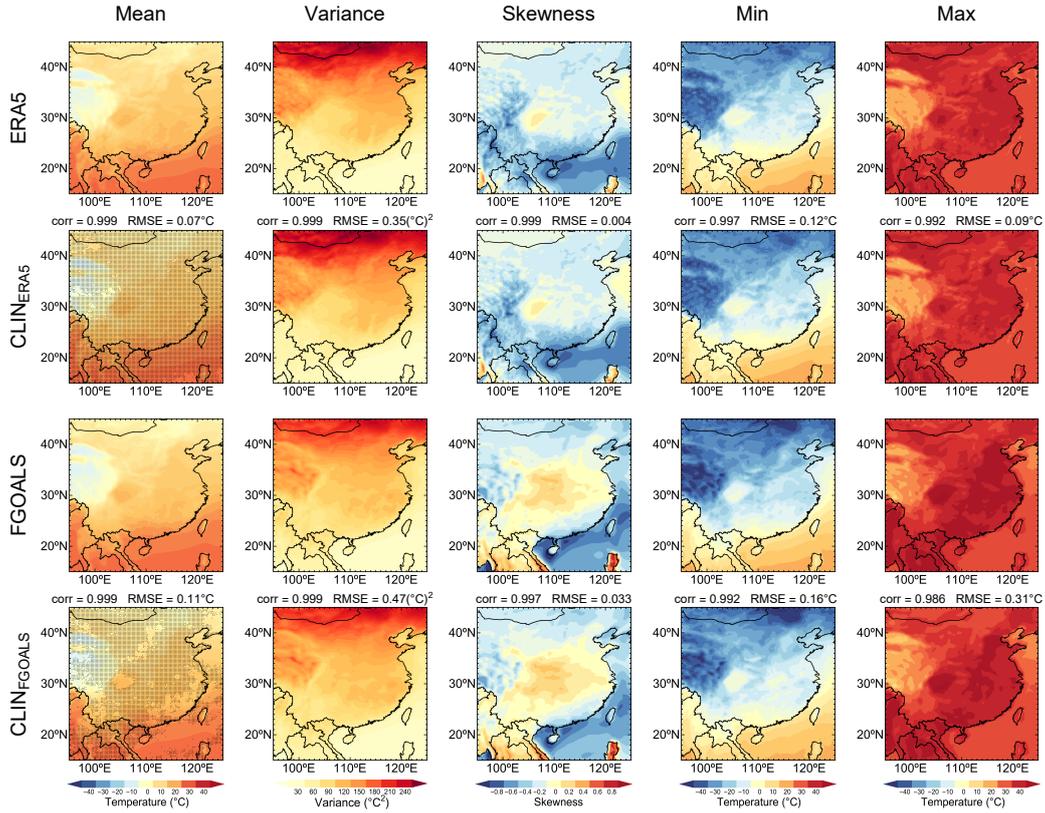
222  $CLIN_{ERA5}$  (Fig. 2 Row 2) and  $CLIN_{FGOALS}$  (Fig. 2 Row 4) can well reproduce the  
 223 considered statistics of their training data, achieving high spatial correlation coefficient  
 224 ( $\sim 0.99$ ) and low root mean squared error ( $\sim 0.1^\circ\text{C}$ ) in matching these statistics. Be-  
 225 sides reproducing the large scale patterns, both models accurately capture high frequency  
 226 local variations influenced by complex topography, such as for mountainous regions and  
 227 coastal areas. Also, the climatology difference between ERA5 and FGOALS are well re-  
 228 produced by the corresponding CLIN models.

229 We further carry out grid-wise Kolmogorov-Smirnov tests to assess whether the gen-  
 230 erated and referential samples are likely to have come from the same underlying distri-  
 231 bution: 96/76% grid points (stippled grids in Fig. 2) within the considered region pass  
 232 a 95% confidence interval test for the  $CLIN_{ERA5}$  and  $CLIN_{FGOALS}$  model. These results  
 233 suggest that the CLIN model can well reproduce climatological distribution of its train-  
 234 ing data at grid scale.

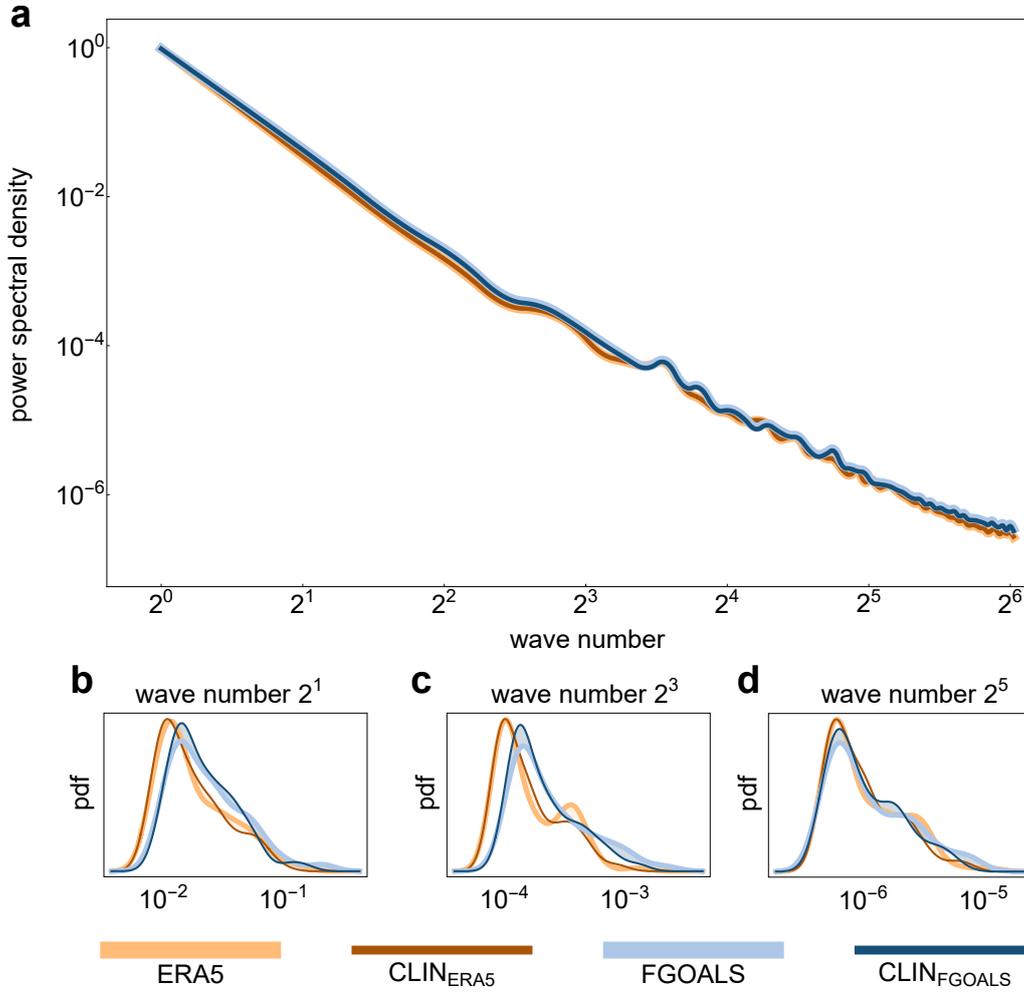
235 We hereafter compare the referential and generated distributions using field-scale  
 236 statistics. We first examine the linear spatial structure of the 2 m temperature spatial  
 237 patterns using a principal component analysis (Supporting Information Fig. S2): we de-  
 238 compose the spatial pattern of the target random variable into a set of orthogonal modes  
 239 that capture the maximum amount of variance, and compare the spatial modes (Em-  
 240 pirical Orthogonal Functions, EOFs), as well as the variance explained by these modes.  
 241 For ERA5, the first to third leading principal components explained 90/2.7/2.0% of the  
 242 total variance. While for  $CLIN_{ERA5}$ , the first to third leading principal components ex-  
 243 plained 91/2.6/1.5% of the total variance, which closely matches results for the ERA5  
 244 referential data. More importantly, we obtain spatial correlation coefficient of 0.994/0.990/0.986  
 245 between the first to third EOF of ERA5 and  $CLIN_{ERA5}$ . While the spatial modes of FGOALS  
 246 differs considerably with ERA5,  $CLIN_{FGOALS}$  closely matches FGOALS: the first to third  
 247 leading principal components explained 83.6/5.1/2.2% or 83.9/4.9/2.1% of the total vari-  
 248 ance for FGOALS or  $CLIN_{FGOALS}$ . The spatial correlation coefficient between the first  
 249 to third EOF of FGOALS and  $CLIN_{FGOALS}$  are 0.999/0.997/0.994. These results sug-  
 250 gest that the CLIN model can well reproduce the linear spatial mode of the considered  
 251 climatological distribution.

252 Lastly, we examine the distribution of spatial variability across different spatial scales  
 253 in the referential/generated dataset: we carry out 2D Fourier transform on the referen-  
 254 tial/generated samples, and draw the radial averaged squared magnitude of the complex  
 255 Fourier coefficients as function of wave numbers (Fig. 3). The radially averaged power  
 256 spectrum density of the considered referential and generated data samples follow a sim-  
 257 ilar power-law scaling, suggesting that the CLIN model can well reproduce the spatial  
 258 variability across scales.

259 To sum up, the analysis of both grid-scale and field-scale statistics demonstrates  
 260 that the CLIN methodology accurately reproduces the essential characteristics and pat-  
 261 terns of the climatological distribution present in the training data. We can thereafter  
 262 leverage this learned climatological prior for the state inference task.



**Figure 2.** Grid-scale comparison of climatological statistics for climate reanalysis (ERA5, Row 1), climate simulation (FGOALS, Row 3), and probabilistic diffusion models trained using these datasets (CLIN<sub>ERA5</sub>, Row 2; and CLIN<sub>FGOALS</sub>, Row 4). The considered statistics are mean, variance, skewness, minimum, and maximum. The spatial correlation coefficient (corr) and root mean squared error (RMSE) between the referential dataset statistics and generated dataset statistics are labeled. Stipples denote grids that pass the Kolmogorov-Smirnov test at 95% confidence interval.



**Figure 3.** Radial averaged power spectrum density as function of wave number for 2 m temperature spatial pattern. **a:** results for ERA5, FGOALS, CLIN<sub>ERA5</sub>, and CLIN<sub>FGOALS</sub> averaged over 100 ensemble members. **b-d:** probability distribution of power spectrum density at wave number 2<sup>1</sup>, 2<sup>3</sup>, 2<sup>5</sup> for ERA5, FGOALS, CLIN<sub>ERA5</sub>, and CLIN<sub>FGOALS</sub>.

263

### 3.2 Inferring weather states using partial observations

264

265

266

267

268

269

270

271

Given a learned climatological prior, we assess how well we can combine it with partial observations to obtain probabilistic estimate of the 2 m temperature spatial patterns. The climatological priors are probabilistic diffusion models trained using climate reanalysis (ERA5) and climate simulation (FGOALS) data. The observations are from 131 operational meteorological stations across China. We randomly select 120 of these stations to inpaint into the generation process, and leave the rest 11 stations for test. For regions without station observations, we consider ERA5 data as benchmark. Below we report case example results (Sec. 3.2.1) and a 1-year round skill assessment (Sec. 3.2.2).

272

#### 3.2.1 Case study

273

274

275

276

277

278

279

280

We consider four case examples covering different hours of a day and different seasons (Fig. 4). To make probabilistic inference of spatial patterns using partial observations, we gradually inpaint station observations into the generation process of  $\text{CLIN}_{\text{ERA5}}$  and  $\text{CLIN}_{\text{FGOALS}}$ , creating 100 ensemble members for each model and each case. We report the ERA5 spatial pattern (Fig. 4 Row 1), the ensemble mean (Fig. 4 Row 2 and 5), the standard deviation of the ensemble (Fig. 4 Row 3 and 6), the mean squared error between ERA5 and the ensemble members (Fig. 4 Row 4 and 7) for  $\text{CLIN}_{\text{ERA5}}$  and  $\text{CLIN}_{\text{FGOALS}}$ .

281

282

283

284

285

286

Both the repainted  $\text{CLIN}_{\text{ERA5}}$  and  $\text{CLIN}_{\text{FGOALS}}$  ensemble mean results closely match the ERA5 spatial pattern, regarding latitudinal gradient, influence of topography, and the land-sea contrast, yielding spatial correlation coefficient of  $0.980 \pm 0.02 / 0.977 \pm 0.02$  for the four considered case examples. These results suggest that the proposed methodology allows effectively propagation of information from limited ( $\sim 1\%$ ) observed locations to a broad range of unobserved parts.

287

288

289

290

291

292

293

294

295

296

297

298

299

300

Next, we test if the CLIN methodology offers reliable uncertainty quantification (Fig. 4 Row 3 and 6). A larger ensemble variance indicates greater uncertainty in the estimate, while a smaller variance suggests more confidence in the estimate. As is expected, geogrids close to observation stations tend to have low ensemble variance, while distant ones may have relatively higher ensemble variance. The information constraint from observations may be blocked by topography, such as for Tibetan Plateau and Tian Shan Mountains. While for plain regions, we can expect a larger extension of observation constraints. We further examine the relationship between the spread of the ensemble members and their estimation skill, by computing the correlation between ensemble variance and ensembles' mean squared error score. The high spread skill correlation for  $\text{CLIN}_{\text{ERA5}}$  ( $0.90 \pm 0.08$ ) and  $\text{CLIN}_{\text{FGOALS}}$  ( $0.94 \pm 0.04$ ) suggest that ensemble spread is a good predictor of model's estimation skill. This means that the CLIN model can capture the underlying uncertainties and provide reliable estimates of spatial estimation confidence.

301

302

303

304

305

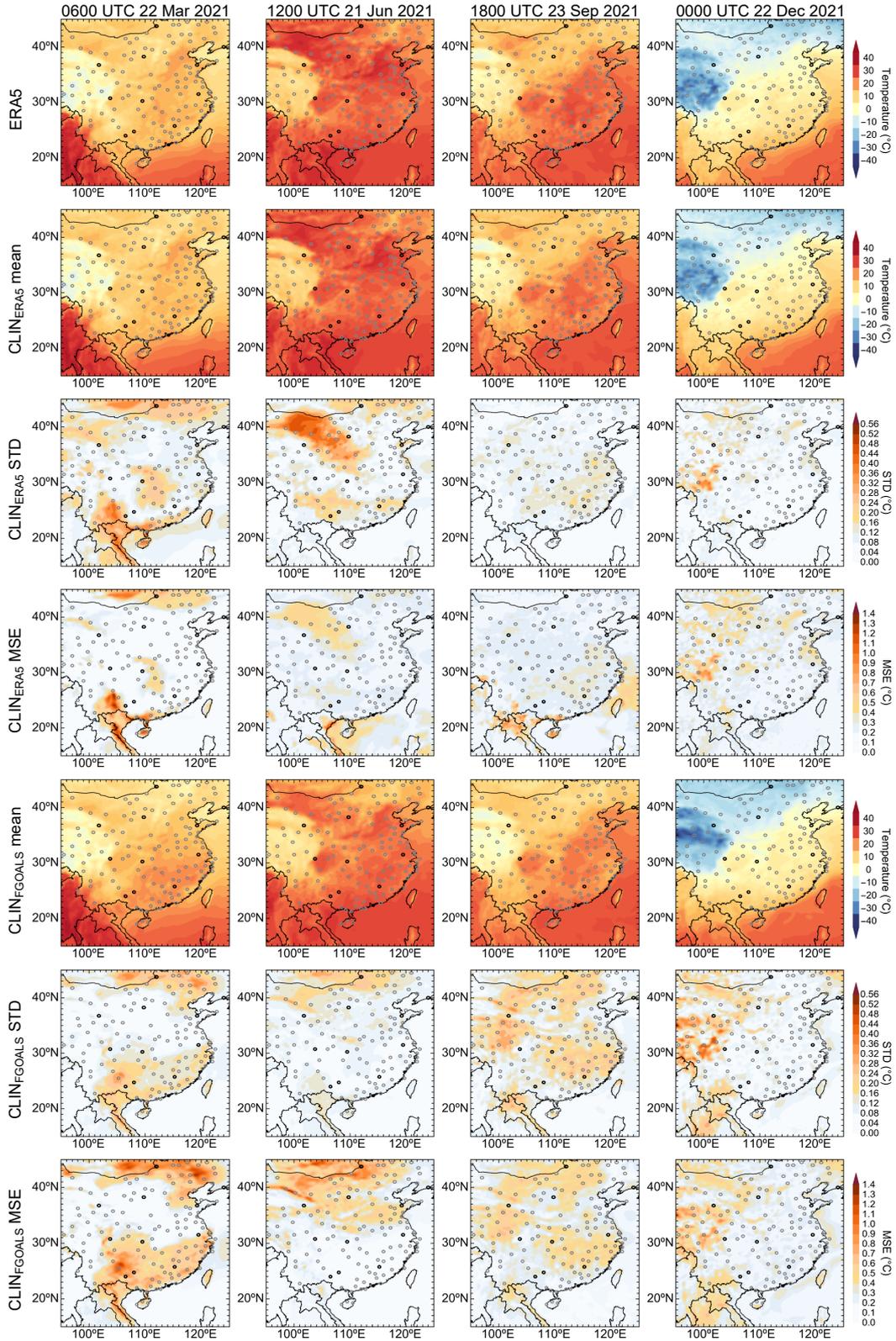
306

307

308

309

To sum up, the case studies confirm that the CLIN methodology can make successful probabilistic inference of 2 m temperature spatial patterns using limited observations. The results are spatially coherent, well-constrained by observations, and offer reliable uncertainty quantification. It is worth noting that there are unneglectable mismatches between station observations and ERA5/FGOALS, regarding either climatological statistics or values. These mismatches introduce domain shift error, which is frequently encountered as we deploy a machine learning model in real-world scenarios where the data distribution differs from the training data. Below we dissect this error source by inpainting with different data sources in a 1-year round evaluation.



**Figure 4.** Case examples for probabilistic inference for 2 m temperature spatial pattern using partial observations. For  $CLIN_{ERA5}$  and  $CLIN_{GOALS}$ , 100 ensemble members are created by repainting observations. The ERA5 spatial pattern (Row 1), the ensemble mean (Row 2 and 5), the standard deviation of the ensemble (Row 3 and 6), the mean squared error between ERA5 and the ensemble members (Row 4 and 7) for  $CLIN_{ERA5}$  and  $CLIN_{GOALS}$  are plotted.

310

### 3.2.2 Skill evaluation

311

312

313

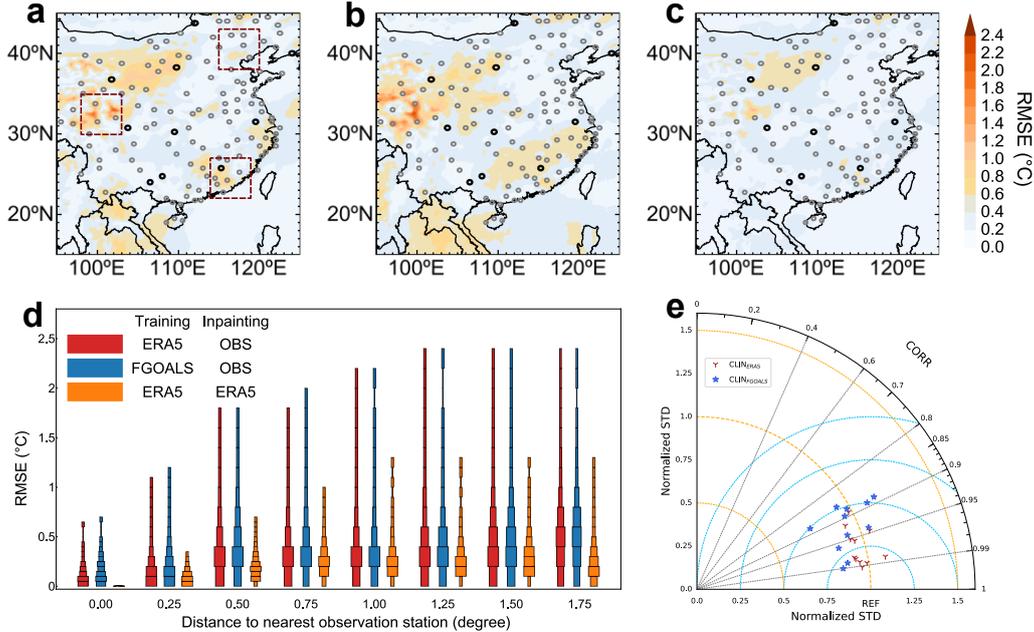
314

315

316

317

We conduct a year-long evaluation of the models' performance in inferring spatial patterns, using data from Year 2021, which are not included in the models' training process. We compare ERA5 with  $CLIN_{ERA5}$  and  $CLIN_{FGOALS}$ , both inpainted using station observations, and present the spatial distribution of their RMSE in Fig. 5a and Fig. 5b. To further investigate different uncertainty sources in the state inference task, we also consider inpainting  $CLIN_{ERA5}$  using ERA5 data at the observation stations. The RMSE between this inpainted  $CLIN_{ERA5}$  and the ERA5 whole-field data is shown in Fig. 5c.



**Figure 5.** Skill evaluation for CLIN models to estimate spatial pattern of 2 m temperature using data for Year 2021. **a:** root mean squared error (RMSE) between ERA5 reanalysis and  $CLIN_{ERA5}$  inpainted using station observations; **b:** RMSE between ERA5 reanalysis and  $CLIN_{FGOALS}$  inpainted using station observations; **c:** RMSE between ERA5 reanalysis and  $CLIN_{ERA5}$  inpainted using ERA5 data at station observations; **d:** distribution of RMSE as function of grid's distance to nearest observation station for the three considered methods; **e:** Taylor diagram comparing the left-out station observations with  $CLIN_{ERA5}$  (orange) and  $CLIN_{FGOALS}$  (blue) results. Both  $CLIN_{ERA5}$  and  $CLIN_{FGOALS}$  are constrained by 120 station observations here. We delineate three representative regions to evaluate the value of observations in Sec. 3.3

318

319

320

321

322

323

324

325

326

The RMSE between ERA5 and observation inpainted  $CLIN_{ERA5}/CLIN_{FGOALS}$  is  $0.25 \pm 0.21^\circ\text{C}/0.31 \pm 0.20^\circ\text{C}$ , suggesting that the CLIN methodology enables accurate spatial pattern estimates. Both models exhibit low uncertainty in plain terrain regions or over the ocean, despite that no ocean observations were applied. This suggests that the learned climatological prior effectively captures the spatial patterns and variability in these regions, allowing the models to make confident estimates using limited and far-away observational constraints. On the other hand, both models exhibit higher uncertainty in regions with complex terrain, such as the Tibetan Plateau and the mountainous areas of Southeast China. Additionally, land areas with complicated terrain but lack

327 ing observational constraints, such as Southeast Asia, also show large uncertainty in the  
 328 model estimates.

329 The uncertainty in state inference comes from the following three sources (Tab. 1).  
 330 The first is domain shift error, which is due to distribution mismatch among data ap-  
 331 plied for model training, data applied for inpainting, and data applied for skill evalua-  
 332 tion. The second is model error, which is due to the approximation/optimization/statistical  
 333 error in applying probabilistic diffusion model to fit climatological prior, or due to er-  
 334 rors in inpainting. These two types of uncertainties are *epistemic*, as they could be re-  
 335 duced by gathering more data, improving the model, or incorporating knowledge about  
 336 data distribution differences. The third source of uncertainty is intrinsic/aleatoric, which  
 337 is due to existence of multiple plausible spatial patterns given partial observational con-  
 338 straints, reflecting the inherent randomness in the system being modeled.

339 To disentangle these uncertainty sources, we consider the following comparisons.

- 340 1. We compare the RMSE of  $\text{CLIN}_{\text{ERA5}}$  (Fig. 5a) and  $\text{CLIN}_{\text{FGOALS}}$  (Fig. 5b).  $\text{CLIN}_{\text{ERA5}}$   
 341 achieves an overall lower RMSE, which can be attributed to a relieved domain shift  
 342 error from the following two aspects: a. compared to FGOALS, ERA5 better matches  
 343 the “true” climatology as partially revealed by the scattered observations; b. we  
 344 consider ERA5 data as “ground truth” for evaluating model performance, which  
 345 gives advantage to CLIN model trained using ERA5 data.
- 346 2. We compare  $\text{CLIN}_{\text{ERA5}}$  inpainted using observation data (Fig. 5a) and  $\text{CLIN}_{\text{ERA5}}$   
 347 inpainted using scattered ERA5 data (Fig. 5c). The latter achieve significantly  
 348 lower RMSE ( $0.19 \pm 0.12^\circ\text{C}$ ), suggesting a relatively low model error and a rel-  
 349 atively low intrinsic uncertainty of the considered task. The difference between  
 350 these cases highlights the domain shift error as the observation distribution dif-  
 351 fers from ERA5.
- 352 3. We compare the performance of  $\text{CLIN}_{\text{ERA5}}$  and  $\text{CLIN}_{\text{FGOALS}}$  in predicting the  
 353 observations at test stations that are excluded during repainting (Fig. 5e). For these  
 354 test stations, both  $\text{CLIN}_{\text{ERA5}}$  and  $\text{CLIN}_{\text{FGOALS}}$  results show high correlation co-  
 355 efficient (0.87-0.99) and low root mean squared error (0.2-0.4 $^\circ\text{C}$ ) with the obser-  
 356 vations, with  $\text{CLIN}_{\text{ERA5}}$  performing slightly better than  $\text{CLIN}_{\text{FGOALS}}$ ;  $\text{CLIN}_{\text{ERA5}}$   
 357 holds a normalized standard deviation close to 1, which closely matches the ob-  
 358 servations, while  $\text{CLIN}_{\text{FGOALS}}$  holds a normalized standard deviation slightly less  
 359 than 1, suggesting a smaller temporal variability.

**Table 1.** Uncertainty sources for state inference using partial observations

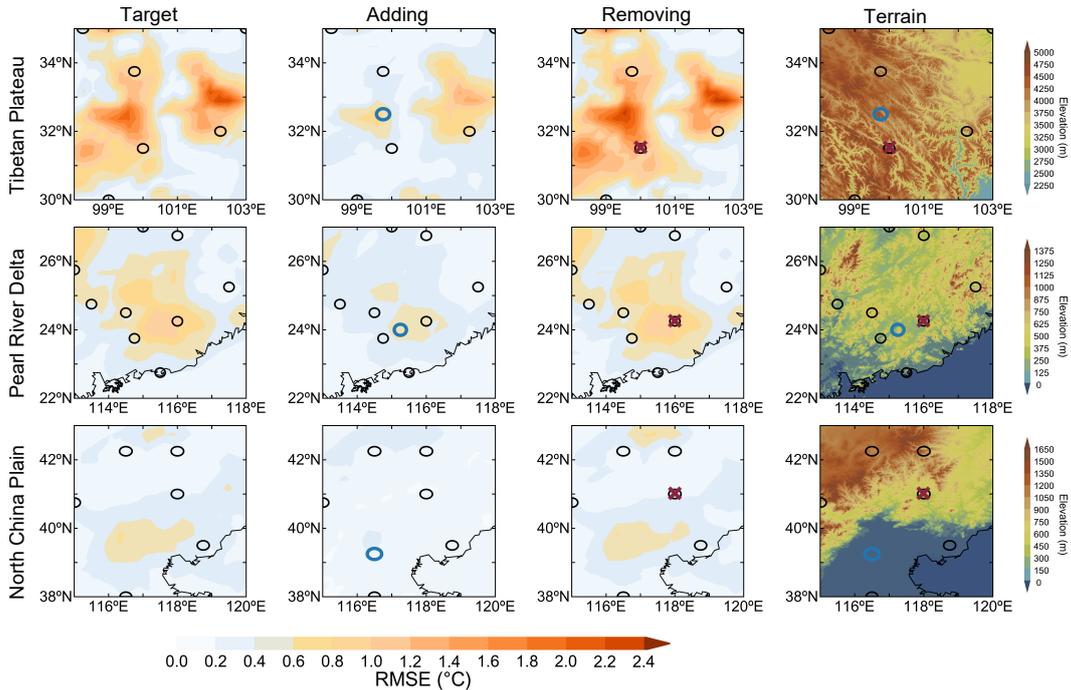
Uncertainty source	Type	Illustration
Domain shift	Epistemic	Distribution mismatch among data applied for model training, data applied for inpainting, and data applied for skill evaluation.
Model error	Epistemic	1. Approximation/optimization/statistical error in fitting climatological prior. 2. Error in constraining the prior with observations.
Intrinsic uncertainty	Aleatoric	Existence of multiple plausible spatial patterns given observational constraints.

360 Finally, we quantify the spatial extension of observational constraints by showing  
 361 models’ RMSE skill as function of grid’s distance to nearest observation station (Fig. 5d).  
 362 We consider  $\text{CLIN}_{\text{ERA5}}$  inpainted using observation data and ERA5 data, as well as  $\text{CLIN}_{\text{FGOALS}}$   
 363 inpainted using observation data. For all these cases, models’ performances at an arbi-  
 364 trary grid depends closely on the grid’s proximity to observations. Meanwhile, there is  
 365 large variation of models’ RMSE skills for grids that are at least  $1^\circ$  away from any ob-  
 366 servation stations. Below we further investigate the value of individual observations in

367 constraining the variability of its nearby spatial patterns, and offer guidelines for bet-  
 368 ter observation planning.

### 369 3.3 On the value of observations

370 We apply the CLIN methodology to quantify the value of observations in constrain-  
 371 ing state estimation uncertainty, using three representative regions delineated in Fig. 5a.  
 372 To achieve this, we add or remove observational stations and evaluate the impact on the  
 373 estimation error (Fig. 6). Here, the first column shows the RMSE spatial pattern for the  
 374 original  $\text{CLIN}_{\text{ERA5}}$  model estimates in each target region; the second column (Adding)  
 375 demonstrates the impact of adding an observation station in a high-error area; the third  
 376 column (Removing) illustrates the effect of removing an existing observation station; the  
 377 fourth column (Terrain) provides a topographical context for each target region.



**Figure 6.** Evaluation of CLIN in reconstructing 2m temperature spatial pattern using different observation setups. Column 1: RMSE between  $\text{CLIN}_{\text{ERA5}}$  inpainted using observation data and ERA5 for three selected regions delineated in Fig. 5. Column 2: RMSE after including a pseudo new observation. This new observation data is from ERA5. Column 3: RMSE after including a pseudo new observation. Column 4: elevation map of the considered regions. The results are based on a year-long (Year 2021) evaluation.

378 For the case of Tibetan Plateau (first row), where the terrain is highly complex,  
 379 with average elevations exceeding 4500 meters, we obtain a relatively high RMSE given  
 380 existing observation constrains, particularly in the central and eastern parts of the re-  
 381 gion. Adding a station in the high-error area significantly reduces the RMSE for a broad  
 382 range of the considered region, this impact is more pronounced here as compared to the  
 383 other two cases, highlighting the importance of observational constraints in areas with  
 384 complex terrain. Removing a station results in a noticeable increase in RMSE in the sur-  
 385 rounding areas. Similarly, the effect of station removal is more evident compared to the

386 other two cases, suggesting that the model heavily relies on the limited observational data  
 387 to constrain its estimates in this complex terrain. The loss of a station in a critical lo-  
 388 cation can greatly impact the model’s ability to capture the local temperature patterns.

389 For the case of Peal River Delta (second row), the terrain is characterized by a mix  
 390 of lowlands and hilly regions, with elevations ranging from 0 to 1000 meters. The orig-  
 391 inal RMSE is low overall, with some higher values in the central and northwest moun-  
 392 tain regions. Adding a station in the high-error area effectively reduces the RMSE. Mean-  
 393 while, removing a station leads to a hardly noticeable increase in RMSE in the surround-  
 394 ing areas.

395 For the case of North China Plain (third row), the northern part is featured by moun-  
 396 tainous terrains exceeding 1000 meters, and the southern part has flat topography and  
 397 homogeneous terrain. Adding a station in the central of southern plain area reduces the  
 398 RMSE significantly, as existing observations are either from the northern mountain ar-  
 399 eas, or is too far away. Same as previous case, removing a station has minimal impact  
 400 on the RMSE distribution.

401 To sum up, we discuss the application of the CLIN methodology to evaluate the  
 402 impact of observational data on state estimation uncertainty across three diverse regions.  
 403 It emphasizes the importance of strategic addition and removal of observational stations  
 404 in improving estimation accuracy, particularly in areas with complex terrain. The find-  
 405 ings highlight how existing observation constraints influence RMSE distribution, with  
 406 significant reductions observed when stations are added in high-error areas. Conversely,  
 407 removal of stations leads to increased RMSE, underscoring the model’s reliance on lim-  
 408 ited observational data. Overall, we provide valuable insights for optimizing the design  
 409 of observation networks, leading to a reduction in uncertainties and biases in weather  
 410 and climate analysis.

## 411 4 Conclusion

412 Accurate state estimation of Earth atmosphere marks a daunting task due to its  
 413 high-dimensionality and chaotic nature. We demonstrated the potential of deep gener-  
 414 ative models, specifically probabilistic diffusion models, in learning the inherent low-dimensional  
 415 statistical structure of atmospheric circulation from climate reanalysis and simulation  
 416 data. By leveraging this learned climatological prior, we developed a methodology named  
 417 CLIN (Climate Inpainting) to effectively infer weather states from partial observations.

418 For the case study of estimating 2 m temperature spatial patterns, the learned cli-  
 419 matological prior accurately reproduced the essential characteristics and patterns of the  
 420 training data at both grid-scale and field-scale. This learned prior effectively captured  
 421 multi-scale climate patterns, providing regularization and stability to the state estima-  
 422 tion task.

423 Combining the learned climatological prior with station observations, CLIN yielded  
 424 strong posterior estimates of 2 m temperature spatial patterns. The estimates were spa-  
 425 tially coherent, well-constrained by observations, and provided reliable uncertainty quan-  
 426 tification. Regions near observation stations exhibited low ensemble variance, indicat-  
 427 ing high confidence in the estimates, while distant regions showed relatively higher en-  
 428 semble variance. The high spread-skill correlation confirmed that the ensemble spread  
 429 was a good predictor of the model’s estimation skill.

430 Moreover, CLIN allowed us to quantify the value of each observation station in re-  
 431 ducing state estimation uncertainty. By adding or removing stations and evaluating the  
 432 impact on the estimation error, we demonstrated the potential of this approach in guid-  
 433 ing the design of optimal observation networks.

Our study showcases the power of deep generative models in extracting and utilizing the information produced by the chaotic evolution of the climate system. The proposed CLIN methodology opens up new opportunities for data-driven weather state estimation, potentially complementing traditional data assimilation approaches.

Future work could focus on extending CLIN to handle indirect observations (i.e., remote sensing) and multiple interdependent variables, incorporating temporal dynamics, and adapting to long-term climate trends. Addressing the computational demands and data requirements of diffusion models is another important direction for making this approach more practical and accessible.

In conclusion, this study demonstrates the immense potential of deep generative models in advancing climate data exploration and tackling complex inference tasks in atmospheric sciences. By learning the intrinsic statistical structure of the climate system, these models can effectively bridge the gap between sparse observations and complete weather state estimates, paving the way for more accurate and efficient climate monitoring and prediction.

## 5 Data Availability

The ERA5 reanalysis data are obtained from the Copernicus Climate Change Service (C3S) Climate Data Store (CDS), accessible at <https://cds.climate.copernicus.eu/>.

The FGOALS model data are obtained from the Coupled Model Intercomparison Project Phase 6 (CMIP6), hosted by the Program for Climate Model Diagnosis and Intercomparison (PCMDI) at Lawrence Livermore National Laboratory (LLNL), accessible at <https://pcmdi.llnl.gov/CMIP6/>.

The observational data are freely available for download from the following website: <http://www.ncdc.noaa.gov/oa/ncdc.html>. The site information used in this study was obtained from the China Meteorological Data Network, hosted by the China National Meteorological Science Data Center (NMDC), accessible at <http://data.cma.cn/>.

## 6 Open Research

Model configuration, analysis scripts, data files used for this study will be publicly available upon accept of the work.

## Acknowledgments

This research is supported by National Key R&D Program of China (Grant NoS. 2023YFC3007700 and 2023YFC3007705). We appreciate the insightful discussions with Dr. Niklas Boers from Technical University of Munich, and with Dr. Bin Wang and Dr. Juanjuan Liu from Chinese Academy of Science.

## References

- Bao, Q., Liu, Y., Wu, G., He, B., Li, J., Wang, L., ... others (2020). Cas fgoals-f3-h and cas fgoals-f3-l outputs for the high-resolution model intercomparison project simulation of cmip6. *Atmospheric and Oceanic Science Letters*, 13(6), 576–581.
- Carrassi, A., Bocquet, M., Bertino, L., & Evensen, G. (2018). Data assimilation in the geosciences: An overview of methods, issues, and perspectives. *Wiley Interdisciplinary Reviews: Climate Change*, 9(5), e535.
- Dhariwal, P., & Nichol, A. (2021). Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34, 8780–8794.

- 478 Eyring, V., Bony, S., Meehl, G. A., Senior, C. A., Stevens, B., Stouffer, R. J., &  
 479 Taylor, K. E. (2016). Overview of the coupled model intercomparison project  
 480 phase 6 (cmip6) experimental design and organization. *Geoscientific Model*  
 481 *Development*, *9*(5), 1937–1958.
- 482 Ghil, M. (2020). Hilbert problems for the climate sciences in the 21st century–20  
 483 years later. *Nonlinear Processes in Geophysics*, *27*(3), 429–451.
- 484 Gilpin, W. (2024). Generative learning for nonlinear dynamics. *Nature Reviews*  
 485 *Physics*, 1–13.
- 486 Hersbach, H., Bell, B., Berrisford, P., Hirahara, S., Horányi, A., Muñoz-Sabater, J.,  
 487 ... others (2020). The era5 global reanalysis. *Quarterly Journal of the Royal*  
 488 *Meteorological Society*, *146*(730), 1999–2049.
- 489 Ho, J., Jain, A., & Abbeel, P. (2020). Denoising diffusion probabilistic models. *Ad-*  
 490 *vances in neural information processing systems*, *33*, 6840–6851.
- 491 Holton, J. R., & Hakim, G. J. (2012). *An introduction to dynamic meteorology*. Aca-  
 492 *ademic press*.
- 493 Kadow, C., Hall, D. M., & Ulbrich, U. (2020). Artificial intelligence reconstructs  
 494 missing climate information. *Nature Geoscience*, *13*(6), 408–413.
- 495 Kanngießer, F., & Fiedler, S. (2024). “seeing” beneath the clouds—machine-  
 496 learning-based reconstruction of north african dust plumes. *AGU Advances*,  
 497 *5*(1), e2023AV001042.
- 498 Kingma, D., Salimans, T., Poole, B., & Ho, J. (2021). Variational diffusion models.  
 499 *Advances in neural information processing systems*, *34*, 21696–21707.
- 500 Law, K., Stuart, A., & Zygalakis, K. (2015). Data assimilation. *Cham, Switzerland:*  
 501 *Springer*, *214*, 52.
- 502 Lugmayr, A., Danelljan, M., Romero, A., Yu, F., Timofte, R., & Van Gool, L.  
 503 (2022). Repaint: Inpainting using denoising diffusion probabilistic models.  
 504 In *Proceedings of the ieee/cvf conference on computer vision and pattern recog-*  
 505 *nition* (pp. 11461–11471).
- 506 Nai, C., Pan, B., Chen, X., Tang, Q., Ni, G., Duan, Q., ... Liu, X. (2024). Reli-  
 507 able precipitation nowcasting using probabilistic diffusion models. *Environmental*  
 508 *Research Letters*.
- 509 Palmer, T. (2017). The primacy of doubt: Evolution of numerical weather prediction  
 510 from determinism to probability. *Journal of Advances in Modeling Earth Sys-*  
 511 *tems*, *9*(2), 730–734.
- 512 Pan, B., Anderson, G. J., Goncalves, A., Lucas, D. D., Bonfils, C. J., Lee, J., ...  
 513 Ma, H.-Y. (2021). Learning to correct climate projection biases. *Journal of*  
 514 *Advances in Modeling Earth Systems*, *13*(10), e2021MS002509.
- 515 Pan, B., Wang, L.-Y., Zhang, F., Duan, Q., Li, X., Pan, X., ... others (2023). Prob-  
 516 abilistic diffusion model for stochastic parameterization—a case example of  
 517 numerical precipitation estimation. *Authorea Preprints*.
- 518 Pan, X., Guo, X., Li, X., Niu, X., Yang, X., Feng, M., ... others (2021). National  
 519 tibetan plateau data center: promoting earth system science on the third pole.  
 520 *Bulletin of the American Meteorological Society*, *102*(11), E2062–E2078.
- 521 Schneider, T., Lan, S., Stuart, A., & Teixeira, J. (2017). Earth system modeling  
 522 2.0: A blueprint for models that learn from observations and targeted high-  
 523 resolution simulations. *Geophysical Research Letters*, *44*(24), 12–396.
- 524 Sohl-Dickstein, J., Weiss, E., Maheswaranathan, N., & Ganguli, S. (2015). Deep  
 525 unsupervised learning using nonequilibrium thermodynamics. In *International*  
 526 *conference on machine learning* (pp. 2256–2265).
- 527 Song, Y., Garg, S., Shi, J., & Ermon, S. (2020). Sliced score matching: A scalable  
 528 approach to density and score estimation. In *Uncertainty in artificial intelli-*  
 529 *gence* (pp. 574–584).
- 530 Song, Y., Sohl-Dickstein, J., Kingma, D. P., Kumar, A., Ermon, S., & Poole, B.  
 531 (2020). Score-based generative modeling through stochastic differential equa-  
 532 tions. *arXiv preprint arXiv:2011.13456*.

- 533 Toth, Z., Talagrand, O., Candille, G., & Zhu, Y. (2003). Probability and ensemble  
534 forecasts. *Forecast verification: A practitioner's guide in atmospheric science*,  
535 *137*, 163.
- 536 Wang, B., Zou, X., & Zhu, J. (2000). Data assimilation and its applications. *Pro-*  
537 *ceedings of the National Academy of Sciences*, *97*(21), 11143–11144.
- 538 Zhang, G., Ji, J., Zhang, Y., Yu, M., Jaakkola, T. S., & Chang, S. (2023). Towards  
539 coherent image inpainting using denoising diffusion implicit models.

# Supporting Information for “Learning to infer weather states using partial observations”

Jie Chao<sup>1,2</sup>, Baoxiang Pan<sup>2</sup>, Quanliang Chen<sup>1</sup>, Shangshang Yang<sup>2,3</sup>, Jingnan Wang<sup>2,4</sup>, Congyi Nai<sup>2,5</sup>, Yue Zheng<sup>6</sup>, Xichen Li<sup>2</sup>, Huiling Yuan<sup>3</sup>, Xi Chen<sup>2</sup>, Bo Lu<sup>7</sup>, Ziniu Xiao<sup>2</sup>

<sup>1</sup>School of Atmospheric Sciences, Chengdu University of Information Technology, Sichuan, China

<sup>2</sup>Institute of Atmospheric Physics, Chinese Academy of Science, Beijing, China

<sup>3</sup>Key Laboratory of Mesoscale Severe Weather, Ministry of Education, and School of Atmospheric Sciences, Nanjing University,

Jiangsu, China

<sup>4</sup>College of Computer, National University of Defense Technology, Hunan, China

<sup>5</sup>Institute of Geographic Sciences and Natural Resources Research, Chinese Academy of Sciences, Beijing, China

<sup>6</sup>Clustertech LTD, Hong Kong, China

<sup>7</sup>National Climate Center, China Meteorological Administration, Beijing, China

## Contents of this file

1. S1 Details of probabilistic diffusion model
2. S2 Parameters of CLIN
3. S3 Model parameter schedule
4. S4 Evaluation metrics
5. Tabel S1 Hyperparameters of Diffusion model

6. Figures S1 Network architecture of diffusion model

7. Figures S2 The first three EOF modes.

## S1. Details of probabilistic diffusion model

Here, we provide detailed mathematical formulations and implementation specifics of the deployed probabilistic diffusion model. For more information and useful learning materials, refer to the works of Sohl-Dickstein et al.(2015), Ho et al. (2020), Song et al. (2020) , Kingma et al. (2021), Ho & Salimans (2022) and Luo (2022).

Diffusion models are probabilistic models that describe the evolution of a stochastic process over time. In the context of deep learning diffusion models, the diffusion process and its reverse process are fundamental concepts.

The diffusion process is the forward process through which a model generates data, typically images, from a simple noise distribution (often Gaussian noise) to the target distribution. A step-by-step derivation is provided below.

First, we define the following Gaussian process to transform the target distribution  $p(\mathbf{x}_0)$  to a prior distribution  $p(\mathbf{x}_T)$ :

$$q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t}\mathbf{x}_{t-1}, \beta_t\mathbf{I}) \quad (1)$$

Here  $\mathbf{x}_{t \in [1, T]}$  are latent variables with increasing noise level;  $\mathcal{N}$  is Gaussian distribution;  $\mathbf{I}$  is identity matrix;  $\beta_t$  is diffusion coefficient, which is pre-defined so that, give large enough  $T$ ,  $p(\mathbf{x}_T|\mathbf{x}_0)$  is drawn close to  $p(\mathbf{x}_T)$ , which is  $\mathbf{x}_0$  agnostic.

$$q(\mathbf{x}_{1:T}|\mathbf{x}_0) = \prod_{t=1}^T q(\mathbf{x}_t|\mathbf{x}_{t-1}) \quad (2)$$

We parameterize the Gaussian encoder with mean  $\mu_t(x_t) = \sqrt{\alpha_t}x_{t-1}$ , and variance  $\Sigma_t(\mathbf{x}_t) = (1 - \alpha_t)\mathbf{I}$ , Here  $\alpha_t = 1 - \beta_t$ . Mathematically, encoder transitions are denoted as:

$$\mathbf{x}_t = \sqrt{\alpha_t}\mathbf{x}_{t-1} + \sqrt{1 - \alpha_t}\boldsymbol{\epsilon}_{t-1} \quad (3)$$

$$= \sqrt{\alpha_t\alpha_{t-1}}\mathbf{x}_{t-2} + \sqrt{1 - \alpha_t\alpha_{t-1}}\bar{\boldsymbol{\epsilon}}_{t-2} \quad (4)$$

$$= \dots \quad (5)$$

$$= \sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}\boldsymbol{\epsilon} \quad (6)$$

$$q(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t}\mathbf{x}_0, (1 - \bar{\alpha}_t)\mathbf{I}) \quad (7)$$

These assumptions depict a systematic process of adding Gaussian noise to the data input over time. As we continue to corrupt the data, it gradually transitions until it is entirely characterized by pure Gaussian noise.

In essence, the reverse process aims to infer the noise distribution that could have generated the observed data. Similar to the diffusion process, the reverse process is represented as:

$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_\theta(\mathbf{x}_t, t), \boldsymbol{\Sigma}_\theta(\mathbf{x}_t, t)) \quad p_\theta(\mathbf{x}_{0:T}) = p(\mathbf{x}_T) \prod_{t=1}^T p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) \quad (8)$$

Here,  $\boldsymbol{\Sigma}_\theta$  is parameterized as an interpolation between its analytical lower and upper bounds (Dhariwal & Nichol, 2021). The optimization of  $\mu_\theta$  involves maximizing the variational lower bound (ELBO) on the log-likelihood of the training samples (Sohl-Dickstein et al., 2015; Kingma et al., 2021).

Then, diffusion model can be optimized by maximizing the ELBO, which can be derived as follows:

$$\log p(\mathbf{x}) = \log \int p(x_{0:T}) dx_{1:T} \quad (9)$$

$$= \log \int \frac{p(x_{0:T})q(x_{1:T}|x_0)}{q(x_{1:T}|x_0)} dx_{1:T} \quad (10)$$

$$= \log \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[ \frac{p(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right] \quad (11)$$

$$\geq \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[ \log \frac{p(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right] \quad (12)$$

$$= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[ \log \frac{p(\mathbf{x}_T)p_\theta(\mathbf{x}_0|\mathbf{x}_1) \prod_{t=2}^T p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_1|\mathbf{x}_0) \prod_{t=2}^T q(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0)} \right] \quad (13)$$

$$= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[ \log \frac{p(\mathbf{x}_T)p_\theta(\mathbf{x}_0|\mathbf{x}_1)}{q(\mathbf{x}_1|\mathbf{x}_0)} + \log \prod_{t=2}^T \frac{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)}{\frac{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)q(\mathbf{x}_t|\mathbf{x}_0)}{q(\mathbf{x}_{t-1}|\mathbf{x}_0)}} \right] \quad (14)$$

$$= \mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} [\log p_0(x_0|\mathbf{x}_1)] + \mathbb{E}_{q(\mathbf{x}_T|\mathbf{x}_0)} \left[ \log \frac{p(\mathbf{x}_T)}{q(\mathbf{x}_T|\mathbf{x}_0)} \right] \\ + \sum_{t=2}^T \mathbb{E}_{q(\mathbf{x}_t, \mathbf{x}_{t-1}|\mathbf{x}_0)} \left[ \log \frac{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)} \right] \quad (15)$$

$$= \mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} [\log p_\theta(x_0|\mathbf{x}_1)] - D_{\text{KL}}(q(\mathbf{x}_T|\mathbf{x}_0) \parallel p(\mathbf{x}_T)) \\ - \sum_{t=2}^T \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t))] \quad (16)$$

We now explain the three terms on the right-hand side of the Eq. 16:

- $\mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} [\log p_\theta(\mathbf{x}_0|\mathbf{x}_1)]$  represents the expected log-likelihood of the initial data  $\mathbf{x}_0$  given the sampled intermediate data  $\mathbf{x}_1$ . For the first step, we have  $\mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} [\log p_\theta(\mathbf{x}_0|\mathbf{x}_1)] = 0$ .

- $D_{\text{KL}}(q(\mathbf{x}_T|\mathbf{x}_0) \parallel p(\mathbf{x}_T))$  denotes the KL divergence between the approximate posterior distribution  $q(\mathbf{x}_T|\mathbf{x}_0)$  and the prior distribution  $p(\mathbf{x}_T)$  at the final time step  $\mathbf{T}$ .

Where  $p(\mathbf{x}_T) \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ , it implies that  $\mathbb{E}_{q(\mathbf{x}_{T-1}|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_T|\mathbf{x}_{T-1}) \parallel p(\mathbf{x}_T))] = 0$

- $\sum_{t=2}^T \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t))]$  represents the sum of the expected KL divergences between the approximate posterior distributions  $q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$  and the

conditional distributions  $p_\theta(\mathbf{x}_t|\mathbf{x}_{t+1})$  for each intermediate time step  $\mathbf{t}$  in the reverse diffusion process.

Given the analysis above, maximizing  $\log p(\mathbf{x})$  can be approximately achieved by minimizing the third term. While minimizing each KL Divergence term individually can be challenging for arbitrary posteriors, we can leverage Bayes' rule to simplify the process:

$$q(x_{t-1}|x_t, x_0) = \frac{q(x_t|x_{t-1}, x_0)q(x_{t-1}|x_0)}{q(x_t|x_0)} \quad (17)$$

$$= \frac{\mathcal{N}(x_t; \sqrt{\alpha_t}x_{t-1}, (1 - \alpha_t)\mathbf{I})\mathcal{N}(x_{t-1}; \sqrt{\alpha_{t-1}}x_0, (1 - \bar{\alpha}_{t-1})\mathbf{I})}{\mathcal{N}(x_t; \sqrt{\alpha_t}x_0, (1 - \bar{\alpha}_t)\mathbf{I})} \quad (18)$$

$$\propto \exp \left\{ - \left[ \frac{(x_t - \sqrt{\alpha_t}x_{t-1})^2}{2(1 - \alpha_t)} + \frac{(x_{t-1} - \sqrt{\alpha_{t-1}}x_0)^2}{2(1 - \bar{\alpha}_{t-1})} - \frac{(x_t - \sqrt{\alpha_t}x_0)^2}{2(1 - \bar{\alpha}_t)} \right] \right\} \quad (19)$$

$$= \exp \left\{ - \frac{1}{2} \left[ \frac{(x_t - \sqrt{\alpha_t}x_{t-1})^2}{1 - \alpha_t} + \frac{(x_{t-1} - \sqrt{\alpha_{t-1}}x_0)^2}{1 - \bar{\alpha}_{t-1}} - \frac{(x_t - \sqrt{\alpha_t}x_0)^2}{1 - \bar{\alpha}_t} \right] \right\} \quad (20)$$

$$= \exp \left\{ - \frac{1}{2} \left( \frac{1}{\frac{(1-\alpha_t)(1-\bar{\alpha}_{t-1})}{1-\bar{\alpha}_t}} \right) \left[ x_{t-1}^2 - 2 \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})x_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)x_0}{1 - \bar{\alpha}_t} x_{t-1} \right] \right\} \quad (21)$$

$$\propto \mathcal{N}(x_{t-1}; \underbrace{\frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})x_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)x_0}{1 - \bar{\alpha}_t}}_{\mu_q(\mathbf{x}_t, \mathbf{x}_0)}, \underbrace{\frac{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t}}_{\Sigma_q(t)}\mathbf{I}) \quad (22)$$

Hence, it is demonstrated that at each step  $\mathbf{x}_{t-1} \sim q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$  follows a normal distribution. We use the KL Divergence between two Gaussian distributions for calculation.

$$\arg \min_{\theta} D_{\text{KL}}(q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) \| p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t)) \quad (23)$$

$$= \arg \min_{\theta} D_{\text{KL}}(\mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_q, \boldsymbol{\Sigma}_q(t)) \| \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_{\theta}, \boldsymbol{\Sigma}_q(t))) \quad (24)$$

$$= \arg \min_{\theta} \frac{1}{2} \left[ \log \frac{|\boldsymbol{\Sigma}_q(t)|}{|\boldsymbol{\Sigma}_q(t)|} - d + \text{tr}(\boldsymbol{\Sigma}_q(t)^{-1} \boldsymbol{\Sigma}_q(t)) + (\boldsymbol{\mu}_{\theta} - \boldsymbol{\mu}_q)^T \boldsymbol{\Sigma}_q(t)^{-1} (\boldsymbol{\mu}_{\theta} - \boldsymbol{\mu}_q) \right] \quad (25)$$

$$= \arg \min_{\theta} \frac{1}{2} \left[ (\boldsymbol{\mu}_{\theta} - \boldsymbol{\mu}_q)^T (\sigma_q^2(t) \mathbf{I})^{-1} (\boldsymbol{\mu}_{\theta} - \boldsymbol{\mu}_q) \right] \quad (26)$$

$$= \arg \min_{\theta} \frac{1}{2\sigma_q^2(t)} \left[ \|\boldsymbol{\mu}_{\theta} - \boldsymbol{\mu}_q\|_2^2 \right] \quad (27)$$

After optimizing the Diffusion Model, the sampling procedure simplifies to sampling Gaussian noise from  $p(\mathbf{x}_T)$  and iteratively running the denoising transitions  $p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t)$  for  $T$  steps to generate a novel  $\mathbf{x}_0$ . In practice, we denote  $\boldsymbol{\mu}_{\theta}$  as function of neural network parameterization for  $\nabla p(\mathbf{x}_t | \mathbf{x}_0)$ , which is commonly known as the *score function* (Y. Song et al., 2020).

## S2. CLIN

In our approach, we merge the acquired climatology prior with station observations to deduce the posterior probability distribution of the target variable. This allows us to jointly modify both observed and unobserved regions throughout the denoising steps, yielding generated samples that are spatially coherent, faithful and adaptive to observation constraints, and uncertainty-aware. The specific parameters of CLIN are presented in the following table. S1.

## S3. Model parameters schedule

We trained the neural network on the NVIDIA Tesla V100 32GB GPU using CUDA version 12.3. The neural network architecture details of the diffusion model are illustrated

in Fig. S1. Typical hyperparameter configurations for diffusion models are often derived from the (Ho & Salimans, 2022).

The specific hyperparameters of the model are presented in the following table. S1.

we embed the time information, and stack the time embedding as an additional channel to all UNet blocks. Each contracting block consists of a long sequence of  $\{C_{3*3} + N + ReLU\}_3$  operations and a short sequence of  $\{C_{1*1}\}_1$  operations, concatenated as a residual block. Here,  $C_{n*n}$  is convolution layer with kernel receptive field of size  $n * n$ .  $N$  is group normalization, ReLU is rectified linear unit function. Each expand block consists of a long sequence of  $\{R_2 + C_{3*3} + N + ReLU\}_3$  operations and a short sequence of  $\{R_2, C_{1*1}\}_1$  operations, concatenated as a residual block. Here,  $R_n$  resizes the data by  $n$  times using linear interpolation. We begin with a channel size of 64 and double/shrink the channel size by 2 along each contracting/expanding block.

## S4. Evaluation metrics

### S4.1 Pearson correlation coefficient (corr)

The Pearson correlation coefficient (*corr*) between prediction  $\hat{x}$  and observation  $x$  is calculated as follows:

$$corr = \frac{\sum_{i=1}^n (\hat{x}_i - \bar{\hat{x}})(x_i - \bar{x})}{\sqrt{\sum_{i=1}^n (\hat{x}_i - \bar{\hat{x}})^2 \cdot \sum_{i=1}^n (x_i - \bar{x})^2}} \quad (28)$$

### S4.2 Root mean square error (RMSE)

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

$$RMSE = \sqrt{\frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n}} \quad (29)$$

### **S4.3 Empirical Orthogonal Function**

Empirical Orthogonal Function (EOF) analysis, also known as Principal Component Analysis (PCA) in some contexts, is a widely used statistical method in various fields, including meteorology, oceanography, climatology, and geophysics.

It is employed to analyze and extract the dominant patterns of variability present in a multivariate dataset, such as spatial patterns in climate data or in oceanographic data. The detailed calculation method for EOF is based on the PrincipalComponents function in Mathematica.

### **S4.4 Kolmogorov-Smirnov test**

The Kolmogorov-Smirnov test (KS test) is a statistical method used to compare the empirical cumulative distribution function (CDF) of a sample dataset with a reference probability distribution or another sample dataset. It is particularly useful for assessing whether the two datasets are drawn from the same underlying distribution or if they differ significantly.

The KS test operates by computing the maximum difference (or maximum deviation) between the two cumulative distribution functions. This maximum difference, often denoted as the KS statistic ( $D$ ), represents the largest vertical distance between the empirical CDF and the theoretical (or reference) CDF. The KS statistic is then compared against critical values from the Kolmogorov-Smirnov distribution, which depends on the sample size and the significance level chosen for the test.

The specific computation method for the Kolmogorov-Smirnov test is derived from Mathematica's `KolmogorovSmirnovTest` function.

### **S4.5 Power spectrum density**

The radial averaged power spectrum density (PSD) is a quantitative measure used in various fields of science and engineering, including signal processing, optics, and geophysics. It provides valuable insights into the distribution of power across different spatial frequencies in a given signal or image. In this paper, the PSD is calculated by first computing the Fourier transform of the signal or image to obtain its frequency domain representation. The power spectrum density is then computed as the squared magnitude of the Fourier transform. The PSD further averages the power spectrum density over concentric circles or spherical shells centered at the origin, hence the term "radial averaged." This averaging process is performed to capture the isotropic characteristics of the signal or image, ensuring that contributions from all directions are considered equally.

The PSD is particularly useful for analyzing signals or images with rotational symmetry or spatial periodicity. By averaging the power spectrum density radially, it becomes possible to discern patterns or structures that are not readily apparent in the original signal or image. Additionally, the PSD can be used to quantify the dominant spatial frequencies present in the signal or image, providing valuable information for further analysis or interpretation.

In short, the radial averaged power spectrum density offers a comprehensive view of the spatial frequency content of a signal or image, facilitating insights into its underlying structure and characteristics.

The specific calculation method for the PSD is derived from the pySTEPS library in Python.

## References

Dhariwal, P., & Nichol, A. (2021). Diffusion models beat gans on image synthesis.

*Advances in neural information processing systems*, 34, 8780–8794.

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*.

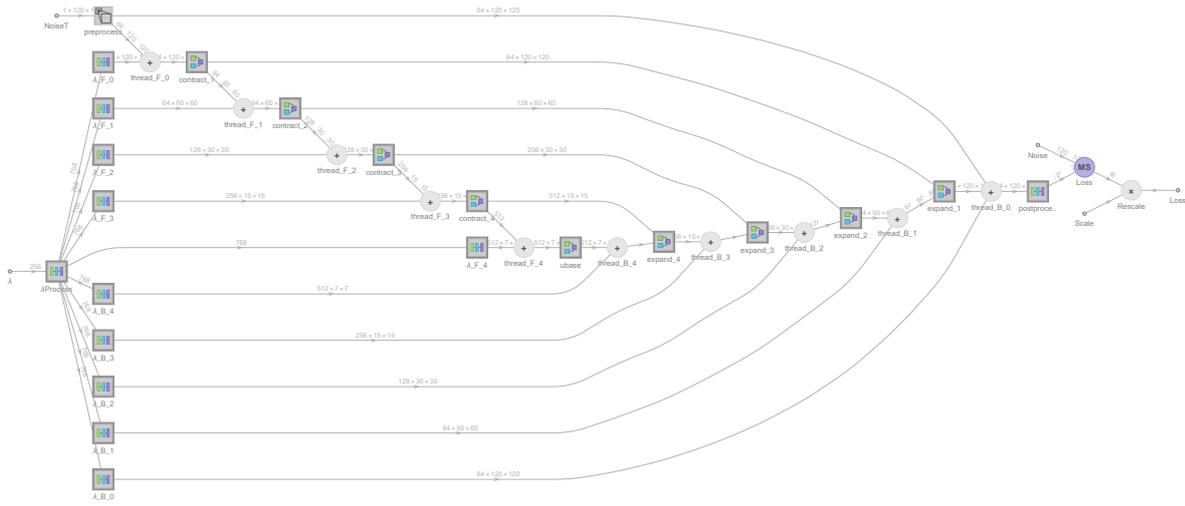
Sohl-Dickstein, J., Weiss, E., Maheswaranathan, N., & Ganguli, S. (2015). Deep unsupervised learning using nonequilibrium thermodynamics. In *International conference on machine learning* (pp. 2256–2265).

Song, J., Meng, C., & Ermon, S. (2020). Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*.

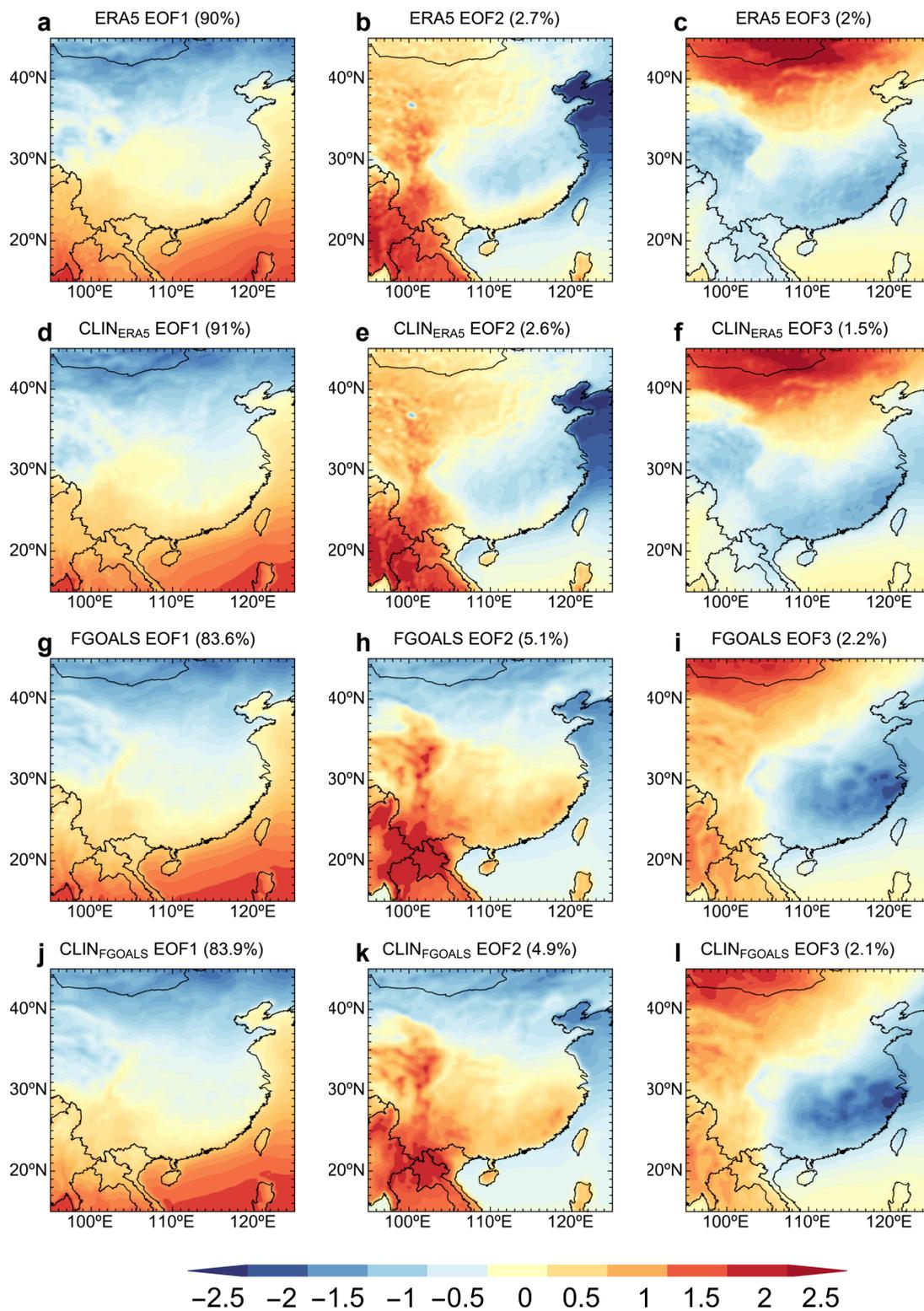
Song, Y., Sohl-Dickstein, J., Kingma, D. P., Kumar, A., Ermon, S., & Poole, B. (2020). Score-based generative modeling through stochastic differential equations. *arXiv preprint arXiv:2011.13456*.

**Table S1.** Hyperparameters of Diffusion model

Hyperparameter	Setting	Parameter
Learning Rate	$10^{-4}$	$\alpha_t^2 = 1 - \sigma_t^2 = \frac{1}{1+e^{-\lambda_t}}$
Batch Size	64	$\lambda_t = -2 \log \tan(at + b)$
Channel	64	$b = \arctan(e^{-\frac{\lambda_{\max}}{2}})$
Optimizer	Adam	$a = \arctan(e^{-\frac{\lambda_{\min}}{2}}) - b$
Number of Iterations	1000	$t = \frac{i}{1000}$ , Where $i = 0, 1, 2, \dots, 1000$
$\lambda_{\min}$	-20	$\text{embedding}(t) = [\sin(2\pi\omega t); \cos(2\pi\omega t)]$
$\lambda_{\max}$	20	$\omega \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$



**Figure S1.** Network Architecture of diffusion model. Each contracting block consists of a long sequence of  $\{C_{3*3} + N + ReLU\}_3$  operations and a short sequence of  $\{C_{1*1}\}_1$  operations, concatenated as a residual block. Here,  $C_{n*n}$  is convolution layer with kernel receptive field of size  $n*n$ .  $N$  is group normalization,  $ReLU$  is rectified linear unit function. Each expand block consists of a long sequence of  $\{R_2 + C_{3*3} + N + ReLU\}_3$  operations and a short sequence of  $\{R_2, C_{1*1}\}_1$  operations, concatenated as a residual block.



**Figure S2.** The first three EOF modes. ERA5 (a-c), *Clin\_ERA5* (d-f), FGOALS (g-i) and CLIN\_FGOALS (j-l).