

# Stochastic Emulators of Spatially Resolved Extreme Temperatures of Earth System Models

Mengze Wang<sup>1</sup>, Andre N. Souza<sup>2</sup>, Raffaele Ferrari<sup>2</sup>, Themistoklis P. Sapsis<sup>1</sup>

<sup>1</sup>Department of Mechanical Engineering, Massachusetts Institute of Technology, Cambridge, MA, USA

<sup>2</sup>Department of Earth, Atmospheric, and Planetary Sciences, Massachusetts Institute of Technology,  
Cambridge, MA, USA

## Key Points:

- Stochastic emulators are developed to estimate the probability distribution of local daily maximum temperature under climate change.
- Coefficients of Empirical Orthogonal Functions are modelled as functions of the global mean temperature, superposed with Gaussian processes.
- Our approach can accurately emulate the quantile anomaly of daily maximum temperature in future scenarios.

---

Corresponding author: Themistoklis P. Sapsis, [sapsis@mit.edu](mailto:sapsis@mit.edu)

**Abstract**

Prediction of extreme events under climate change is challenging but essential for risk management of natural disasters. Although earth system models (ESMs) are arguably our best tool to predict climate extremes, their high computational cost restricts the application to project only a few future scenarios. Emulators, or reduced-complexity models, serve as a complement to ESMs that achieve a fast prediction of the local response to various climate change scenarios. Here we propose a data-driven framework to emulate the full statistics of spatially resolved climate extremes. The variable of interest is the near-surface daily maximum temperature. The spatial patterns of temperature variations are assumed to be independent of time and extracted using Empirical Orthogonal Functions (EOFs). The time dependence is encoded through the coefficients of leading EOFs which are decomposed into long-term seasonal variations and daily fluctuations. The former are assumed to be functions of the global mean temperature, while the latter are modelled as Gaussian stochastic processes with temporal correlation conditioned on the season. The emulator is trained and tested using the simulation data in CMIP6. By generating multiple realizations, the emulator shows significant performance in predicting the temporal evolution of the probability distribution of local daily maximum temperature. Furthermore, the uncertainty of the emulated statistics is quantified to account for the internal variability. The emulation accuracy in testing scenarios remains consistent with the training datasets. The performance of the emulator suggests that the proposed framework can be generalized to other climate extremes and more complicated scenarios of climate change.

**Plain Language Summary**

Extreme events in the global climate system, such as heat waves and hurricanes, cause incalculable losses every year. Conventional climate models, called Earth System Models (ESMs), are our best tools to predict how climate change may affect the occurrence rate of extreme events in the future. However, these models are relatively slow and expensive to run. We present a framework to design emulators, or reduced-complexity models, to efficiently predict the complete statistics of climate extremes on spatially-resolved grids. Once trained using a few simulations generated from ESMs, the emulator can be used to predict climate change scenarios that were not included in the training data. Our approach is demonstrated for near-surface daily maximum temperature data. The mean, variance, and extreme values of the temperature generated by the emulator are very similar to the statistics generated by ESMs. Furthermore, the emulator provides a speedy quantification of the uncertainty of the predicted statistics. The performance of the emulator suggests that our framework can be generalized to other types of extreme events in the climate system.

**1 Introduction**

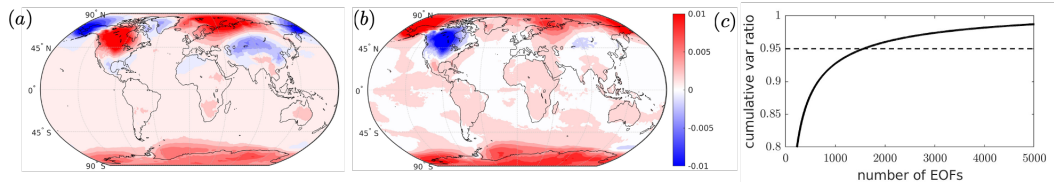
Unprecedented climate extremes, associated with anthropogenic global warming, have been observed worldwide, such as the Russian heatwaves in 2010 and the record-breaking Atlantic hurricane season in 2020 (Meehl & Tebaldi, 2004; Barriopedro et al., 2011; Reed et al., 2022). The annual losses from such weather- and climate-related disasters have surged dramatically, escalating from several billion dollars in 1980 to 200 billion in 2020 (Allen et al., 2012; AON, 2020), not to mention the incalculable loss of lives. Effectively managing the risks of extreme events and minimizing their associated damages necessitates accurate quantification of their likelihood in a rapidly changing global climate. Despite the increased frequency of extreme weather events, their probability at a given time and location is still very low, and thus quantifying their risks requires large ensembles of numerical simulations for very long time horizons. The need for ensembles of simulations amplifies the already high computational cost associated with running full-

64 scale Earth System Models (ESMs) and restricts their application to a limited number  
65 of climate change scenarios. In contrast, emulators, or reduced-complexity models, pro-  
66 vide a more efficient evaluation of the statistics of extreme events in response to more  
67 diverse scenarios. In the present work, we develop a multivariate Gaussian stochastic em-  
68 ulator that estimates the probability distribution of local daily maximum temperature  
69 on spatially-resolved grids.

70 Climate emulators can be broadly categorized by the spatial resolution of their pro-  
71 jections. The first type of emulators, also known as simple climate models (SCMs), fo-  
72 cus on modelling how global or regional mean fields are influenced by the concentrations  
73 of greenhouse gases, emissions of aerosols, and natural effective radiative forcing vari-  
74 ations (Meinshausen et al., 2011; Seneviratne et al., 2016). A majority of these emula-  
75 tors have been systematically compared in the Reduced Complexity Model Intercompar-  
76 ison Project (Z. R. Nicholls et al., 2020; Z. Nicholls et al., 2021), by evaluating their pre-  
77 diction accuracy of the global mean temperature. Based on this type of emulators, in-  
78 teractive models have been developed for policymakers and stakeholders to actively ex-  
79 amine the impact of energy, economic and public policies on climate change (Kapmeier  
80 et al., 2021; Rooney-Varga et al., 2021).

81 The second type of emulators specialize in predicting the response of local variables  
82 to climate change. The most widely used method for this type of emulator is pattern scal-  
83 ing, where the climate variables at different locations are assumed as independent lin-  
84 ear functions of the global mean temperature (Mitchell, 2003). Therefore, the global mean  
85 temperature predicted by the first-type emulators can be used as an input for pattern  
86 scaling, facilitating localized climate predictions in response to a variety of emission sce-  
87 narios. Over time, the framework of pattern scaling has evolved to encompass a broader  
88 range of techniques. These advances include the adoption of response functions to ac-  
89 count for past trajectories of CO<sub>2</sub> (Castruccio et al., 2014; Freese et al., 2024), the use  
90 of Matern covariance functions for modeling spatial correlation (Alexeeff et al., 2018),  
91 and the incorporation of internal variability through autoregressive processes or the spec-  
92 trum of principal components analysis (Beusch et al., 2020; Link et al., 2019). As mod-  
93 ern machine learning methods emerge, researchers have explored diverse architectures  
94 to enhance the accuracy of local climate emulation, utilizing inputs ranging from globally-  
95 averaged emissions to spatial distribution of aerosols. Most of these machine learning  
96 models have been evaluated on the benchmark datasets, with ClimateBench (Watson-  
97 Parris et al., 2022) and ClimateSet (Kaltenborn et al., 2023) being the most commonly  
98 used ones. Compared with pattern scaling, neural networks can provide a more accu-  
99 rate emulation of certain variables, such as the global precipitation, when trained on suf-  
100 ficiently large ensembles of simulations (Lütjens et al., 2024) albeit with a compromise  
101 in the model complexity.

102 Both classes of emulators have been typically used to predict time-averaged quan-  
103 tities. Only a few recent studies have explored emulating the statistics of climate extremes,  
104 such as the annual maximum temperature and the duration of hot waves within a year  
105 (Tebaldi et al., 2020; Quilcaille et al., 2022). Furthermore, no prior work has been re-  
106 ported on the emulation of probability distribution of local climate variables, which con-  
107 stitutes the primary objective of our research. We introduce a stochastic model to em-  
108 ulate the statistics of climate extremes, utilizing temperature-related extreme events as  
109 a prototypical application. We first extract the empirical orthogonal functions (EOF)  
110 (Lorenz, 1956; Hannachi et al., 2007) of the spatial patterns of near-surface daily max-  
111 imum temperature (TMX) fields to reduce the dimensionality of the system while main-  
112 taining a high spatial resolution. Driven by the observed nearly-Gaussian character of  
113 the EOF statistics (conditioned over season and year), we model the temporal evolution  
114 of the EOF coefficients as Gaussian stochastic processes (Mohamad & Sapsis, 2015; Arbabi  
115 & Sapsis, 2022), characterized by long-term trends, seasonal variations, and colored noise.  
116 The mean, variance and covariance of the EOF coefficients are parameterized using the



**Figure 1.** (a,b) The first and second spatial EOFs of daily maximum temperature, computed using CNRM-CM6-1-HR simulation data. (c) Cumulative variance ratio represented by leading EOFs.

117 global mean temperature and season, thus generalizing our emulator to more diverse cli-  
 118 mate change scenarios. A similar framework has been applied to emulate monthly-averaged  
 119 temperature and humidity (Geogdzhayev et al., 2024). Our work will focus on daily max-  
 120 imum temperature and its full statistics.

121 The content of this paper is organized as follows. In §2 we introduce the simula-  
 122 tion data used for training and testing the emulator. The mathematical framework of  
 123 the emulator is described in §3, including the dimensionality reduction method in §3.1  
 124 and stochastic modeling of time series in §3.2. The emulation results are presented in  
 125 §4, followed by a summary of the main conclusions and discussion in 5.

## 126 2 Data

127 Among all the ESMs in Coupled Model Intercomparison Project Phase 6 (CMIP6),  
 128 we adopted the CNRM-CM6-1-HR and MPI-ESM1-2-LR model outputs as our reference  
 129 dataset. Both models achieved reasonable skill scores on simulating the statistics of cli-  
 130 mate extremes according to a recent evaluation of the performance of CMIP6 models (Wehner  
 131 et al., 2020). The CNRM-CM6-1-HR model provides the highest spatial resolution (nom-  
 132 inal resolution 50km) among CMIP6 models, which best fits our needs to develop a spatially-  
 133 resolved emulator. However, this model only has one realization available, which is in-  
 134 sufficient to assess the influence of climate internal variability on the emulator. The MPI-  
 135 ESM1-2-LR data feature a large ensemble of realizations, although the spatial resolu-  
 136 tion (250km) is problematic for studying climate extremes. Therefore, the majority of  
 137 our results will focus on emulation of CNRM-CM6-1-HR data, while the large ensemble  
 138 data of MPI-ESM1-2-LR will be utilized to investigate the impact of internal vari-  
 139 ability and ensemble size on the performance of the emulator.

140 Two variables are collected from the CNRM-CM6-1-HR and MPI-ESM1-2-LR model  
 141 outputs: (i) Near-surface daily mean temperature (the *tas* variable in CMIP6), used to  
 142 compute the global mean temperature; (ii) Near-surface daily maximum temperature (the  
 143 *tasmax* variable in CMIP6). Here “near surface” refers to two-meter height according  
 144 to the CMIP6 convention. The CMIP6 simulations cover a historical period from 1850  
 145 to 2014, followed by a set of future scenarios until 2100. The CNRM-CM6-1-HR model  
 146 offers only one realization for both the historical period and each future scenario, whereas  
 147 the MPI-ESM1-2-LR model provides 50 realizations. To train the emulator, we utilize  
 148 the simulation data within the historical period and the SSP5-8.5 future scenario for each  
 149 ESM. The SSP1-2.6 future scenario is utilized for testing purposes.

### 3 Methods

#### 3.1 Data pre-processing: dimensionality reduction

Since we focus on the near-surface temperature, the spatial location  $\mathbf{x}$  is described by the latitude and longitude coordinates,  $\mathbf{x} = (\theta, \varphi)$ , where  $\theta \in [-\pi/2, \pi/2]$  and  $\varphi \in [0, 2\pi)$ . The time step size is one day, and the number of days since 01/01/1850 0:00 is represented as  $t$ . The daily maximum temperature (TMX) at location  $\mathbf{x}$  and time  $t$  for the ensemble member  $\omega$  is denoted as  $q(\mathbf{x}, t, \omega)$ . The climatological mean  $\bar{q}(\mathbf{x}, t)$  is extracted by phase-averaging TMX for the same calendar day and location across the historical period, 1850-2014, and over the entire ensemble. In other words, at an arbitrary time  $t$ ,  $\bar{q}(\mathbf{x}, t) = \bar{q}(\mathbf{x}, \text{mod}(t, 365))$ . The fluctuations of TMX are decomposed as superposition of Empirical Orthogonal Functions (EOFs),  $\phi_i(\mathbf{x})$ ,

$$q'(\mathbf{x}, t, \omega) := q(\mathbf{x}, t, \omega) - \bar{q}(\mathbf{x}, t) = \sum_i a_i(t, \omega) \phi_i(\mathbf{x}). \quad (1)$$

In order to compute the EOFs, we construct the spatial covariance function  $\mathcal{R}(\mathbf{x}, \mathbf{x}^*)$  that quantifies the covariance between fluctuating TMX at two arbitrary locations  $\mathbf{x}$  and  $\mathbf{x}^*$ ,

$$\mathcal{R}(\mathbf{x}, \mathbf{x}^*) = \langle q'(\mathbf{x}, t, \omega) q'(\mathbf{x}^*, t, \omega) \rangle_{t\omega}. \quad (2)$$

The notation  $\langle \cdot \rangle_{t\omega}$  represents averaging over time and the ensemble. The EOFs are defined as the eigenfunctions of  $\mathcal{R}(\mathbf{x}, \mathbf{x}^*)$ , taking into account the curvature of the Earth's surface  $S$ ,

$$\int_S \mathcal{R}(\mathbf{x}, \mathbf{x}^*) \phi_i(\mathbf{x}^*) \cos \theta^* d\theta^* d\varphi^* = \lambda_i \phi_i(\mathbf{x}). \quad (3)$$

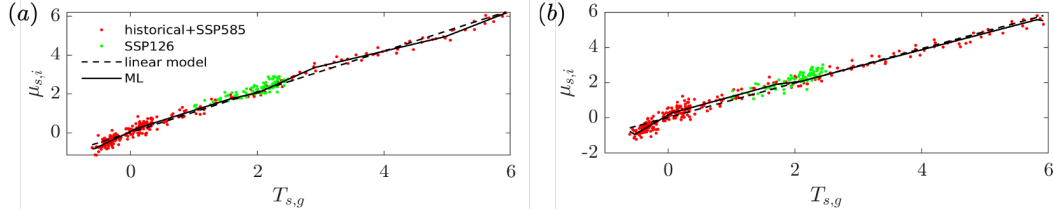
The coefficient of each EOF at time  $t$  is obtained by projecting  $q'(\mathbf{x}, t, \omega)$  onto  $\phi_i(\mathbf{x})$ ,

$$a_i(t, \omega) = \int_S q'(\mathbf{x}, t, \omega) \phi_i(\mathbf{x}) \cos \theta d\theta d\varphi. \quad (4)$$

Similar to the climatological mean, the EOFs are also computed from the historical data. However, we only utilize the TMX snapshots on every five days, rather than daily data, because TMX on adjacent days are highly correlated. Our choice of five-day interval is based on the observation that on this timescale the autocorrelation coefficient of TMX at most locations decreases to approximately 0.5 (Kalvová & Nemesšová, 1998), striking a reasonable balance between data independence and comprehensive representation of temperature variability.

For CNRM-CM6-1-HR data, since only one realization is available, the number of snapshots ( $1.2 \times 10^4$ ) is much smaller than the number of grids ( $2.6 \times 10^5$ ). As such, it is unnecessary to store the large covariance matrix (2), and the method of snapshots is adopted to solve the eigenvalue problem (3) more efficiently. Specifically, we compute the temporal covariance matrix of  $q'$ , whose size is the square of the number of snapshots. The eigen-decomposition of the temporal covariance matrix is then performed to get its eigenvalues and eigenfunctions, which can be linearly transformed to get the eigenpairs  $(\lambda_i, \phi_i(\mathbf{x}))$  of the spatial covariance  $\mathcal{R}(\mathbf{x}, \mathbf{x}^*)$ . More details can be found in Sirovich (1987) and Taira et al. (2020). For MPI-ESM1-2-LR data, the number of grids ( $1.8 \times 10^4$ ) is comparable or smaller than the total number of snapshots ( $1.2 \times 10^4 \times$  the number of realizations adopted), and we directly solve equation (3) to obtain the eigenfunctions of the spatial covariance.

The first two EOFs of the CNRM-CM6-1-HR data are visualized in figure 1(a,b). They account for 2.9% and 2.7% of the total variance, respectively. Both EOFs are reminiscent of the Arctic Oscillation/Northern Hemisphere Annular Mode (Thompson & Wallace, 1998) and the Southern Hemisphere Annular Mode (Fogt & Marshall, 2020). Unlike previous studies that focused on the first few EOFs to extract the physically significant modes (Wallace & Gutzler, 1981; Amaya, 2019), our objective is to reconstruct the



**Figure 2.** Jun-Aug mean of (a) the first and (b) second EOF coefficients in each year of CNRM-CM6-1-HR dataset, from 1850 to 2100, plotted versus the global mean temperature. Red dots: true seasonal mean obtained from the historical and SSP5-8.5 scenario. Green dots: SSP1-2.6 scenario. Black dashed line: linear regression; Solid line: machine-learned function.

193 full probability distribution of local TMX with sufficient accuracy and efficiency. There-  
 194 fore, we retain the first 2,000 EOFs for the CNRM-CM6-1-HR model, which altogether  
 195 represent approximately 95% of the total variance (figure 1c) of the respective datasets.

### 196 3.2 Multivariate Gaussian stochastic emulator of EOF time series

197 Assuming the climatological mean and EOFs remain invariant with respect to time  
 198 and future scenarios, our stochastic emulator of the daily maximum temperature is for-  
 199 mulated as,

$$\hat{q}(\mathbf{x}, t, \hat{\omega}) = \bar{q}(\mathbf{x}, t) + \sum_{i=1}^I \hat{a}_i(t, \hat{\omega}) \phi_i(\mathbf{x}). \quad (5)$$

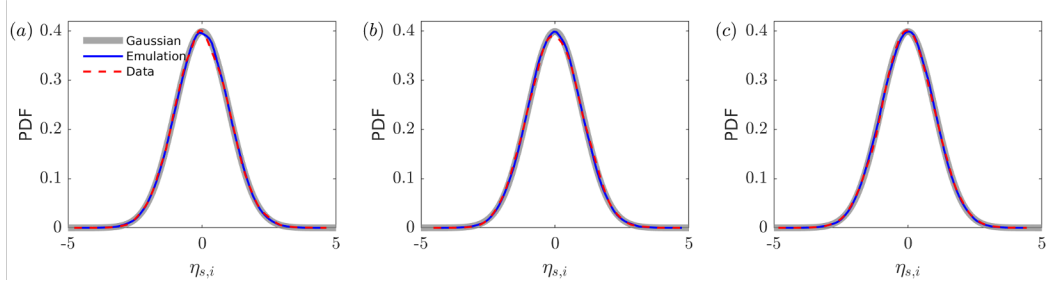
200 A notable difference between equation (5) and the decomposition of true TMX fluctu-  
 201 ations (1) is the EOF coefficient, where  $a$  is the true coefficient obtained from projec-  
 202 tion (4) and  $\hat{a}$  is estimated from the emulator. The emulation index  $\hat{\omega}$  is also different  
 203 from the ensemble member  $\omega$ , since the emulator can be used to generate more realiza-  
 204 tions than the training data.

205 The time series of  $\hat{a}$  in season  $s$  and for a given global mean temperature, is mod-  
 206 elled as superposition of long-term trends and Gaussian-distributed daily fluctuations  
 207 that encode temporal correlation:

$$\hat{a}_{s,i}(t, \hat{\omega}) = \hat{\mu}_{s,i}(T_{s,g}) + \hat{\sigma}_{s,i}(T_{s,g}) \sum_{j=1}^I \hat{l}_{s,ij} \hat{\eta}_{s,j}(t, \hat{\omega}), \quad i = 1, 2, \dots, I. \quad (6)$$

208 The subscript  $s = 1, 2, 3, 4$  corresponds to Northern Hemisphere spring (Mar-May), sum-  
 209 mer (Jun-Aug), autumn (Sep-Nov), and winter (Dec-Feb) respectively. The seasonal mean  
 210  $\hat{\mu}_{s,i}$  and variance  $\hat{\sigma}_{s,i}^2$  are parameterized as a function of the seasonally-averaged global  
 211 mean temperature,  $T_{s,g}$ . The correlation between the  $i$ th and  $j$ th EOFs in season  $s$  is  
 212 assumed constant and accounted for by  $\hat{l}_{s,ij}$ . The daily fluctuations of the EOF coeffi-  
 213 cients are modelled as superposition of Gaussian autoregressive processes  $\hat{\eta}_{s,j}(t, \hat{\omega})$ . Here  
 214  $\hat{\eta}_{s,j}$  and  $\hat{\eta}_{s,k}$  are uncorrelated when  $j \neq k$ , and the time series of  $\hat{\eta}_{s,j}$  are emulated us-  
 215 ing the autocorrelation computed from training data. Specifically, consider a time win-  
 216 drow in season  $s$  of the  $y$ -th year, denoted as  $t \in [t_{ys}, t_{ys} + N_s]$ . The starting time,  $t_{ys}$ ,  
 217 corresponds to the first day of each season: Mar 1st, Jun 1st, Sep 1st, and Dec 1st, for  
 218  $s = \{1, 2, 3, 4\}$ . The duration of each time window,  $N_s$ , is given by  $N_s = \{92, 92, 91, 90\}$   
 219 days respectively. Within  $t \in [t_{ys}, t_{ys} + N_s]$ , the emulated daily fluctuations  $\hat{\eta}_{s,j}(t, \hat{\omega})$   
 220 satisfy

$$\hat{\eta}_{s,j}(t, \hat{\omega}) = \sum_{n=1}^{t-t_{ys}} c_{s,j}(n) \hat{\eta}_{s,j}(t-n, \hat{\omega}) + g_{s,j}(n) \epsilon_{s,j}(n), \quad \epsilon_{s,j}(n) \sim \mathcal{N}(0, 1), \quad t \in [t_{ys}, t_{ys} + N_s]. \quad (7)$$



**Figure 3.** Probability density function (PDF) of the 1st, 2nd, and 500th component of the Jun-Aug  $\eta_s$ : (a)  $\eta_{2,1}$ , (b)  $\eta_{2,2}$ , (c)  $\eta_{2,500}$ . Red dashed lines: PDF computed using CNRM-CM6-1-HR historical and SSP5-8.5 future scenario data, from 1850 to 2100; gray lines: Gaussian fit of  $\eta_{2,i}$  data; blue lines: PDF of 10 emulations of 1850-2100  $\hat{\eta}_{2,i}$

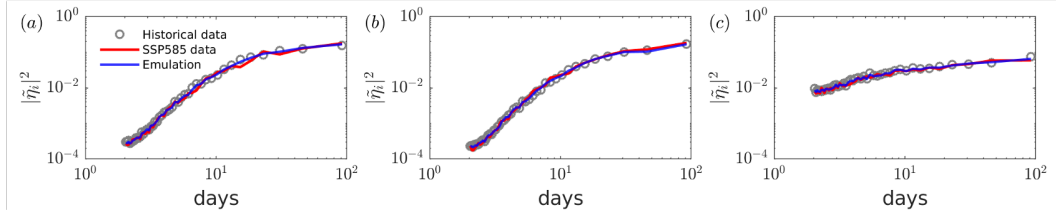
221 The parameters  $c_{s,j}(n)$  and  $g_{s,j}(n)$  are independent of the year and will be estimated  
 222 from the training data, while the standard normal random number  $\epsilon_{s,j}(n)$  varies with  
 223 the year and the emulation. The emulator (6) can also be written more compactly in vec-  
 224 tor form,

$$\hat{\mathbf{a}}_s(t, \hat{\omega}) = \hat{\boldsymbol{\mu}}_s(T_{s,g}) + \hat{\mathbf{D}}_s(T_{s,g}) \hat{\mathbf{L}}_s \hat{\boldsymbol{\eta}}_s(t, \hat{\omega}), \quad (8)$$

225 where  $\hat{\mathbf{a}}_s$ ,  $\hat{\boldsymbol{\mu}}_s$ , and  $\hat{\boldsymbol{\eta}}_s$  are  $I \times 1$  column vectors. The notation  $\hat{\mathbf{D}}_s$  is a diagonal matrix,  
 226 and each element on the diagonal is  $\hat{\sigma}_{s,i}$ . The matrix  $\hat{\mathbf{L}}_s$  is lower triangular, where each  
 227 entry corresponds to  $\hat{l}_{s,ij}$ .

228 It is important to emphasize here that the formulated emulator is conditionally Gaus-  
 229 sian, i.e. for a fixed season and global mean temperature, the daily fluctuations are, by  
 230 design, normally distributed. While this does not necessarily imply that long term statis-  
 231 tics will have a Gaussian character, since we also have the variation of the global mean  
 232 temperature, it does not allow for the possibility of daily temperature extremes that have  
 233 (for a given season and global mean temperature) a non-Gaussian distribution, e.g. fol-  
 234 low heavy tails. For the present context, direct comparisons suggest that this is an accept-  
 235 able assumption. However, for other variables this aspect may introduce limitations. We  
 236 plan to extend the framework to address these potential limitations in future work.

237 The unknown parameters (which are functions of  $T_{s,g}$ ) in the emulator (6,7) are  
 238 estimated using the true EOF coefficients  $a_i(t, \omega)$  (4) and the global mean temperature  
 239  $T_{s,g}$  from 1850 to 2100 (historical and SSP5-8.5 scenario). Given  $a_i(t, \omega)$  data, we first  
 240 compute the actual seasonal mean  $\mu_{s,i}$  and standard deviation  $\sigma_{s,i}$  in each year, aver-  
 241 aged over the entire ensemble. Two examples of the Jun-Aug mean  $\mu_{s,i}$  versus the cor-  
 242 responding  $T_{s,g}$  are shown in figure 2 (red dots). These relationships are mostly linear  
 243 and independent of the future scenario (SSP1-2.6 shown in green dots), which motivate  
 244 us to regress  $\hat{\mu}_{s,i}$  as a linear function of  $T_{s,g}$  (black dashed lines). Similar linear rela-  
 245 tionships are also observed for the variance  $\sigma_{s,i}^2$  and also for higher-ranked EOFs. Nonlin-  
 246 ear functions are also attempted using fully-connected neural networks. For each  $\hat{\mu}_{s,i}$  or  
 247  $\hat{\sigma}_{s,i}^2$ , the neural network is designed with two hidden layers, each containing three neu-  
 248 rons, utilizing the ReLU activation function. The learned nonlinear functions are shown  
 249 as black solid lines in figure 2, which provide slightly better agreement with the train-  
 250 ing data. A more systematic comparison of the emulation results using linear and non-  
 251 linear functions will be provided in §4.1. We also explored alternative network architec-  
 252 tures with varying numbers of layers and neurons, as well as different activation func-  
 253 tions, including Sigmoid and Tanh. However, these modifications did not yield signif-  
 254 icant improvements and the associated results are not shown.



**Figure 4.** Spectra of the 1st, 2nd, and 500th component of the Jun-Aug  $\eta_s$ : (a)  $\eta_{2,1}$ , (b)  $\eta_{2,2}$ , (c)  $\eta_{2,500}$ . Gray circles: spectra averaged using CNRM-CM6-1-HR historical (1850-2014) data; red lines: CNRM-CM6-1-HR SSP5-8.5 (2015-2100) data; blue lines: 10 emulations of 2015-2100 spectra.

255 After extracting the variation of the seasonal mean and standard deviation in re-  
 256 sponse to the global mean temperature,  $\hat{\mu}_{s,i}(T_{s,g})$  or  $\hat{\sigma}_{s,i}^2(T_{s,g})$ , we remove these trends  
 257 from the true EOF coefficients, resulting in the residuals  $(a_{s,i} - \hat{\mu}_{s,i}) / \hat{\sigma}_{s,i}$ . We then evalu-  
 258 ate their cross-correlations,

$$\hat{\Sigma}_s = \left\langle \hat{\mathbf{D}}_s^{-1} (\mathbf{a}_s - \hat{\boldsymbol{\mu}}_s) (\mathbf{a}_s - \hat{\boldsymbol{\mu}}_s)^\top \hat{\mathbf{D}}_s^{-\top} \right\rangle_{t\omega}, \quad \hat{\Sigma}_s = \hat{\mathbf{L}}_s \hat{\mathbf{L}}_s^\top, \quad (9)$$

259 The time average is performed from 1850 to 2100 for each season respectively. While the  
 260 actual cross correlations fluctuate over time, they remain statistically stationary for most  
 261 EOFs, justifying the choice of a constant matrix model. Generalization of (9) to time-  
 262 dependent correlations requires large-ensemble data and will be discussed in §4.2. The  
 263 last equality in (9) is a Cholesky decomposition of  $\hat{\Sigma}_s$ . Multiplying the residuals by  $\hat{\mathbf{L}}_s^{-1}$   
 264 produces uncorrelated time series,

$$\boldsymbol{\eta}_s(t, \omega) = \hat{\mathbf{L}}_s^{-1} \hat{\mathbf{D}}_s^{-1} (\mathbf{a}_s(t, \omega) - \hat{\boldsymbol{\mu}}_s), \quad (10)$$

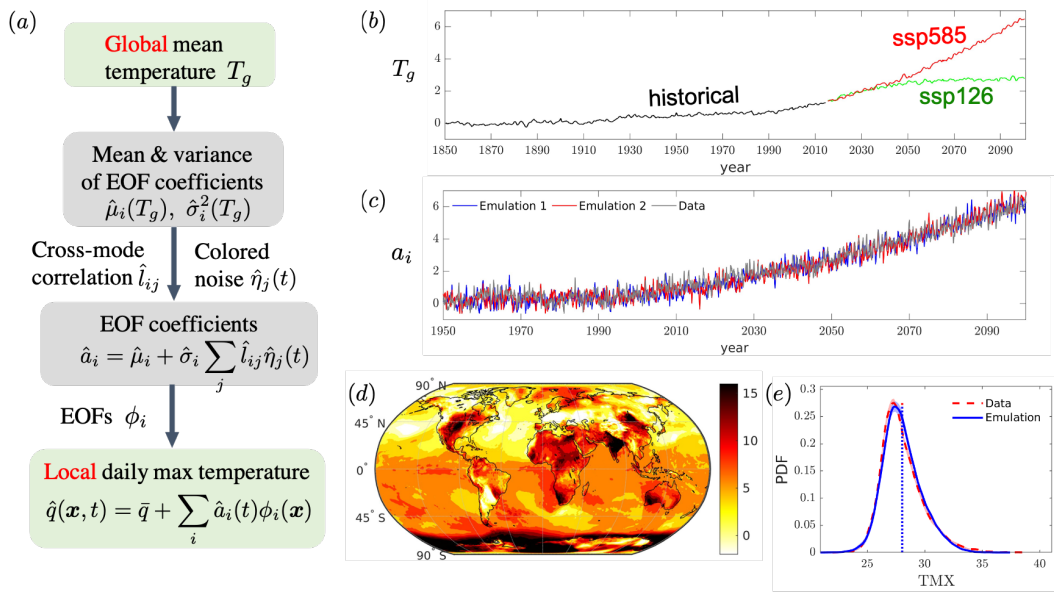
265 which satisfies

$$\langle \boldsymbol{\eta}_s(t, \omega) \boldsymbol{\eta}_s(t, \omega)^\top \rangle_{t\omega} = \mathbf{I}. \quad (11)$$

266 Here  $\mathbf{I}$  is an identity matrix with a size equal to the number of adopted EOFs. In other  
 267 words, each entry of  $\boldsymbol{\eta}_s(t, \omega)$  has unit variance, and different entries are uncorrelated.

268 To justify our assumption that  $\eta_{s,j}(t, \omega)$  in season  $s$  can be modelled as Gaussian  
 269 processes (equation 7) with the same autocorrelations across different years, we evalu-  
 270 ate the statistics  $\eta_{s,j}(t, \omega)$  in figure 3,4. The probability density functions of the 1st, 2nd,  
 271 and 500th component of Jun-Aug  $\eta_{s,j}$  are computed using historical and SSP5-8.5 scen-  
 272 ario data, from 1850 to 2100. The profiles are plotted by red dashed lines in figure 3,  
 273 which almost overlap with the fitted Gaussian distributions (gray lines). While not shown  
 274 here, the other components of  $\eta_{s,j}(t, \omega)$  also exhibit approximately Gaussian distribu-  
 275 tions. To examine the time dependence of the second-order statistics of each component  
 276 of  $\boldsymbol{\eta}_s$ , we compute the Fourier spectra of  $\boldsymbol{\eta}_s$  in Jun-Aug of each year and average them  
 277 over two distinct time windows, 1850-2014 and 2015-2100 of SSP5-8.5 scenario. As vi-  
 278 sualized in figure 4, the spectra of three components of  $\boldsymbol{\eta}_s$  remain approximately unchanged  
 279 over time. Therefore, the statistics averaged over the entire 1850-2100 period are used  
 280 to generate the surrogate Gaussian processes  $\hat{\eta}_{s,j}$  that represent stochastic realiza-  
 281 tions of daily fluctuations. Simulation of the Gaussian processes is based on the exact  
 282 time-domain method which utilizes the autocorrelation of  $\boldsymbol{\eta}_s$ . This approach has been  
 283 demonstrated more robust against uncertainty of statistics than the frequency-domain  
 284 method (Percival, 1993). The PDFs of the simulated  $\hat{\boldsymbol{\eta}}_s$  in figure 3 (blue lines) indeed  
 285 follow Gaussian distribution, and the Fourier spectra of the simulated processes align  
 286 with the true spectra, as illustrated in figure 4.





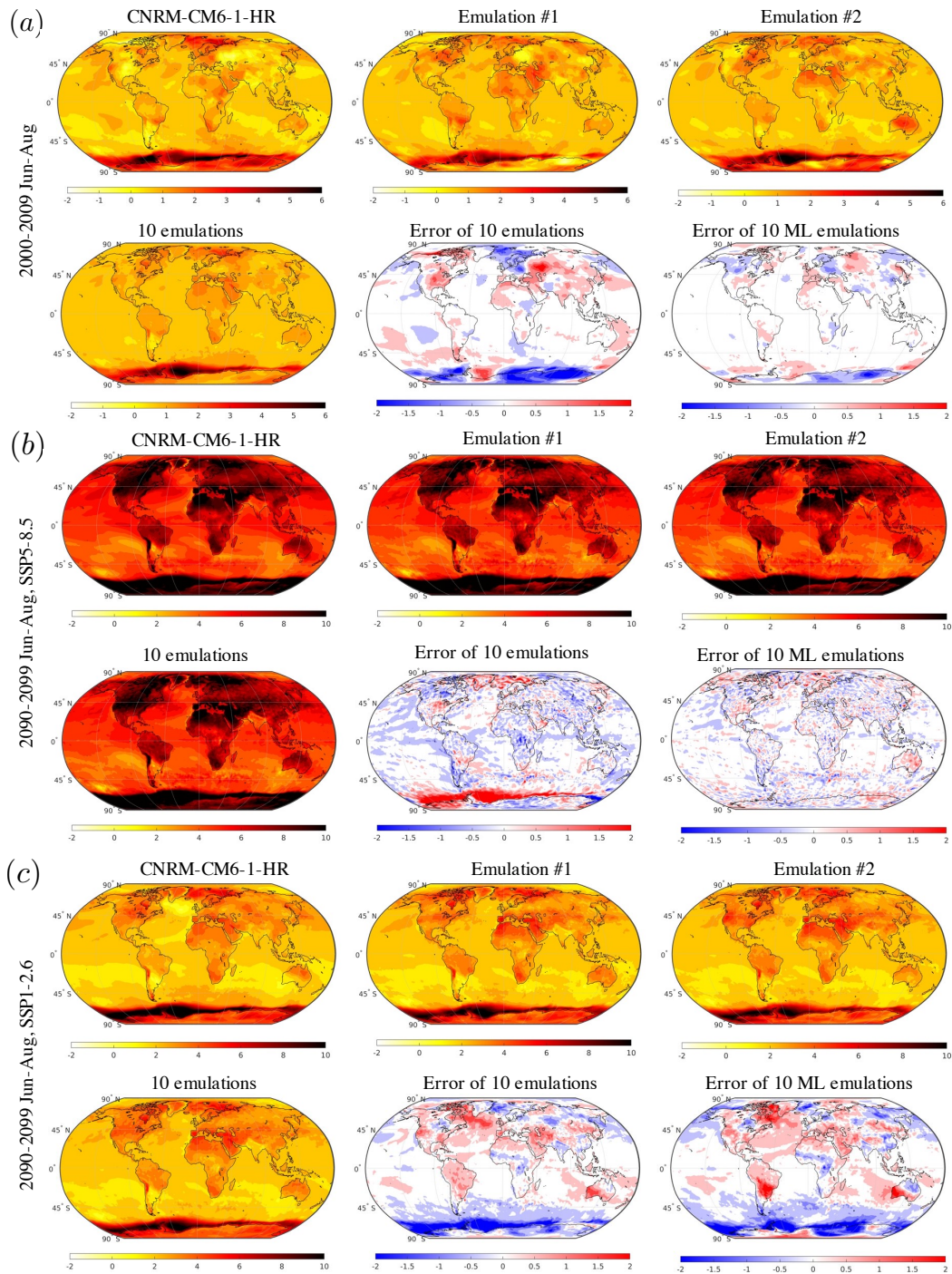
**Figure 5.** (a) Flow chart showing the structure of the emulator. Given the global mean temperature  $T_g$ , the emulator predicts the local daily maximum temperature on spatially-resolved grids. (b) One-year moving average of the global mean temperature, shown for different scenarios. (c) Example time series of the true and emulated EOF coefficients. (d) Sample outputs from the emulator: reconstruction of the TMX field. (e) An example of the probability density function of local TMX, averaged in Jun-Aug over a ten-year window. The vertical lines mark the mean values.

287 The steps of the emulation are summarized schematically in figure 5a. Starting from  
 288 the temporal evolution of the global mean temperature (panel b), the seasonal mean and  
 289 variance of the EOF coefficients are estimated from the learned relationships  $\hat{\mu}_{s,i}(T_{s,g})$ ,  
 290  $\hat{\sigma}_{s,i}^2(T_{s,g})$ . The daily fluctuations are constructed as the stochastic autoregressive pro-  
 291 cesses  $\hat{\eta}_{s,j}(t, \omega)$ , which are scaled by  $\hat{l}_{s,ij}$  and superposed to account for the cross cor-  
 292 relation between different EOFs. Combining the scaled daily fluctuations with long-term  
 293 trends, we obtain the emulated time series of the EOF coefficients, exhibiting the same  
 294 first and second order statistics as the true time series (panel c). Given the time series  
 295 and shape of EOFs, the final output of the emulator is the temporal evolution of grid-  
 296 ded local TMX. A sample snapshot of TMX is visualized in panel d. To acquire converged  
 297 probability distribution of local TMX, especially for the tails that represent extreme events,  
 298 the statistics are computed by averaging over a decadal window in time and a  $1^\circ \times 1^\circ$   
 299 region in space. Panel e shows a sample comparison between the emulated and true prob-  
 300 ability density function (PDF). The blue region marks the uncertainty of the distribu-  
 301 tion, estimated by performing multiple emulations. We note the non-Gaussian charac-  
 302 ter of the target and approximated PDF, which is the result of considering the statis-  
 303 tics over a time window that the global average temperature changes.

## 304 4 Emulation results

### 305 4.1 Emulation of CNRM-CM6-1-HR dataset

306 The performance of the emulator is firstly evaluated in detail for Jun-Aug, when  
 307 TMX is the most extreme in Northern Hemisphere. Results in other seasons will be briefly  
 308 discussed at the end of this section. To differentiate between the emulator that adopts

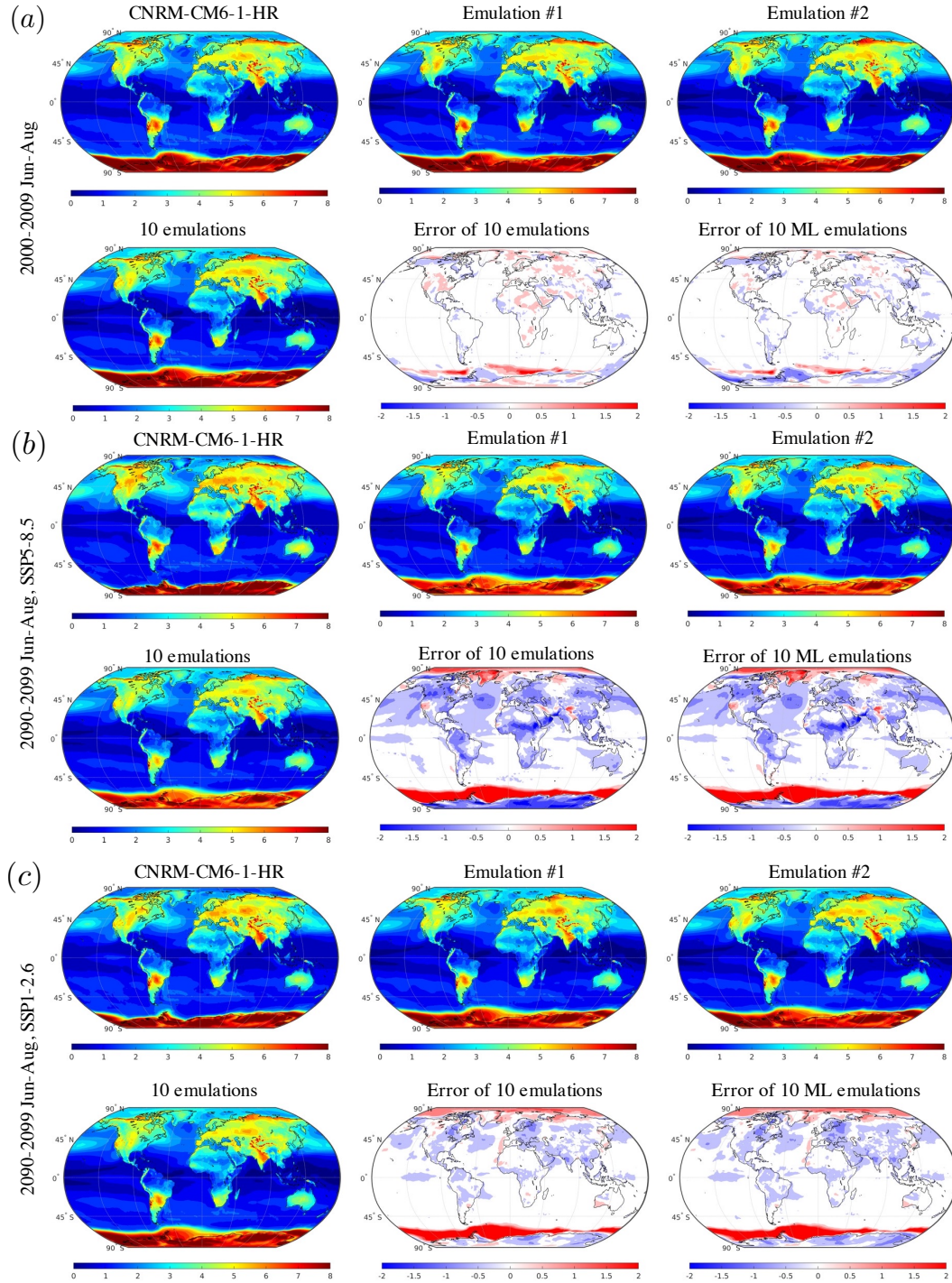


**Figure 6.** Mean anomaly of Jun-Aug daily maximum temperature, averaged over (a) 2000-2009, (b) 2090-2099 of the SSP5-8.5 scenario, and (c) 2090-2099 of the SSP1-2.6 scenario. Each subfigure shows the true mean from CNRM-CM6-1-HR ESM, two sample emulations, average of 10 emulations, error of 10 emulations, and the error of 10 ML emulations. Reference: 1850-1900 Jun-Aug mean TMX.

309 linear and nonlinear model for the long-term trends, the former is referred to as “em-  
 310 ulation” and the latter is denoted as “machine-learning (ML) emulation”. Figure 6 shows  
 311 the mean of local TMX across three decadal periods: 2000-2009 within the historical pe-  
 312 riod, 2090-2099 of the SSP5-8.5 scenario, and 2090-2099 of the SSP1-2.6 scenario. The  
 313 reference mean, computed from the CNRM-CM6-1-HR data, is compared against two  
 314 sample emulations, the average of ten emulations, and ML emulations. The emulator ac-  
 315 curately captures the evolution of local TMX under both high and low warming scenar-  
 316 ios. Significant anomalies in regions such as the Arctic, western coast of South Amer-  
 317 ica, North Africa, West Asia and Southern Ocean are well reproduced. Errors are within  
 318 1°C at most locations, with the highest errors reaching 2°C. Using ML model for the sea-  
 319 sonal mean and variance appreciably improves the emulation accuracy. Despite train-  
 320 ing on historical and SSP5-8.5 data only, the emulator performance on the unseen SSP1-  
 321 2.6 scenario demonstrates its potential for application across various climate change path-  
 322 ways.

323 The errors of the emulated mean in figure 6(a-c) arise from different contributions.  
 324 In figure 6a, the discrepancy between the emulations and the true mean mainly origi-  
 325 nates from the modeling assumption that the seasonal mean is fully determined by the  
 326 global mean temperature,  $\hat{\mu}_{s,i}(T_{s,g})$ . As discussed in §3.2 (c.f. figure 2), a single global  
 327 mean temperature  $T_{s,g}$  can correspond to multiple values of the mean EOF coefficients  
 328  $\mu_{s,i}$ , due to the internal variability of the climate system and the neglected influence of  
 329 the past global mean temperature or emission history. The internal variability of the CNRM-  
 330 CM6-1-HR simulation is difficult to quantify, since only one realization is available. How-  
 331 ever, the variability captured by the emulator can be readily assessed by performing mul-  
 332 tiple emulations. Comparing the pattern of errors with the two emulations in figure 6a,  
 333 we observe that most high-error regions also exhibit high variability, such as Europe and  
 334 the Southern Ocean. In addition, the error magnitude aligns with the variability, indicat-  
 335 ing that the error can be further reduced if more realizations of the ESM are avail-  
 336 able for training the emulator and computing the local statistics. In figure 6b, smaller-  
 337 scale fluctuation of the errors become more apparent, which stems from the changing shape  
 338 of the leading EOFs under different warming conditions. Recall that the EOFs were com-  
 339 puted only using the historical data. The leading historical EOFs adopted in the emu-  
 340 lator may represent a lower variance in the SSP5-8.5 scenario, which results in higher  
 341 emulation errors contributed by truncating EOFs. This issue can be mitigated by includ-  
 342 ing SSP5-8.5 data into the calculation of EOFs, though similar errors might recur when  
 343 the emulator is applied to unseen scenarios. The error in SSP1-2.6 scenario (figure 6c)  
 344 is slightly higher than SSP5-8.5, due to the trained model of long-term trends not be-  
 345 ing optimal for SSP1-2.6. The error of ML emulations are even higher than linear emu-  
 346 lations for SSP1-2.6, such as in South America, which indicates that the superior per-  
 347 formance of ML emulator in SSP5-8.5 is likely due to overfitting. Nevertheless, the sen-  
 348 sitivity of the seasonal mean to warming condition is modest, and the emulation error  
 349 remains the same order of magnitude across different scenarios.

350 The standard deviation of local TMX is presented in figure 7. In historical peri-  
 351 ods, such as 2000-2009 shown in figure 7a, the standard deviation is reconstructed accu-  
 352 rately for most locations. The error from ten emulations is almost identical to the ML  
 353 emulations, suggesting a predominantly linear relationship between the variance of most  
 354 EOF coefficients and the global mean temperature,  $\hat{\sigma}_{s,i}(T_{s,g})$ . From 2000-2009 to 2090-  
 355 2099 in SSP5-8.5 scenario (panel b), the standard deviation slightly increases in most re-  
 356 gions, such as North America, North Africa and West Asia. In contrast, the standard  
 357 deviation in Greenland and Southern Ocean shows a significant reduction, likely due to  
 358 diminished ice coverage (Räisänen, 2002; Gao et al., 2015). These trends are consistent  
 359 with the observational data (Huntingford et al., 2013) and ESM simulations using other  
 360 models (Olonscheck & Notz, 2017). The performance of the emulator is the least sat-  
 361 isfactory in regions associated with the most significant trends. For example, the enhanced  
 362 variance in North Africa is not captured, and the decreasing trend in the Southern Ocean



**Figure 7.** Standard deviation of Jun-Aug daily maximum temperature, averaged over (a) 2000-2009, (b) 2090-2099 of the SSP5-8.5 scenario, and (c) 2090-2099 of the SSP1-2.6 scenario. Each subfigure shows the true mean from CNRM-CM6-1-HR ESM, two sample emulations, average of 10 emulations, error of 10 emulations, and the error of 10 ML emulations.

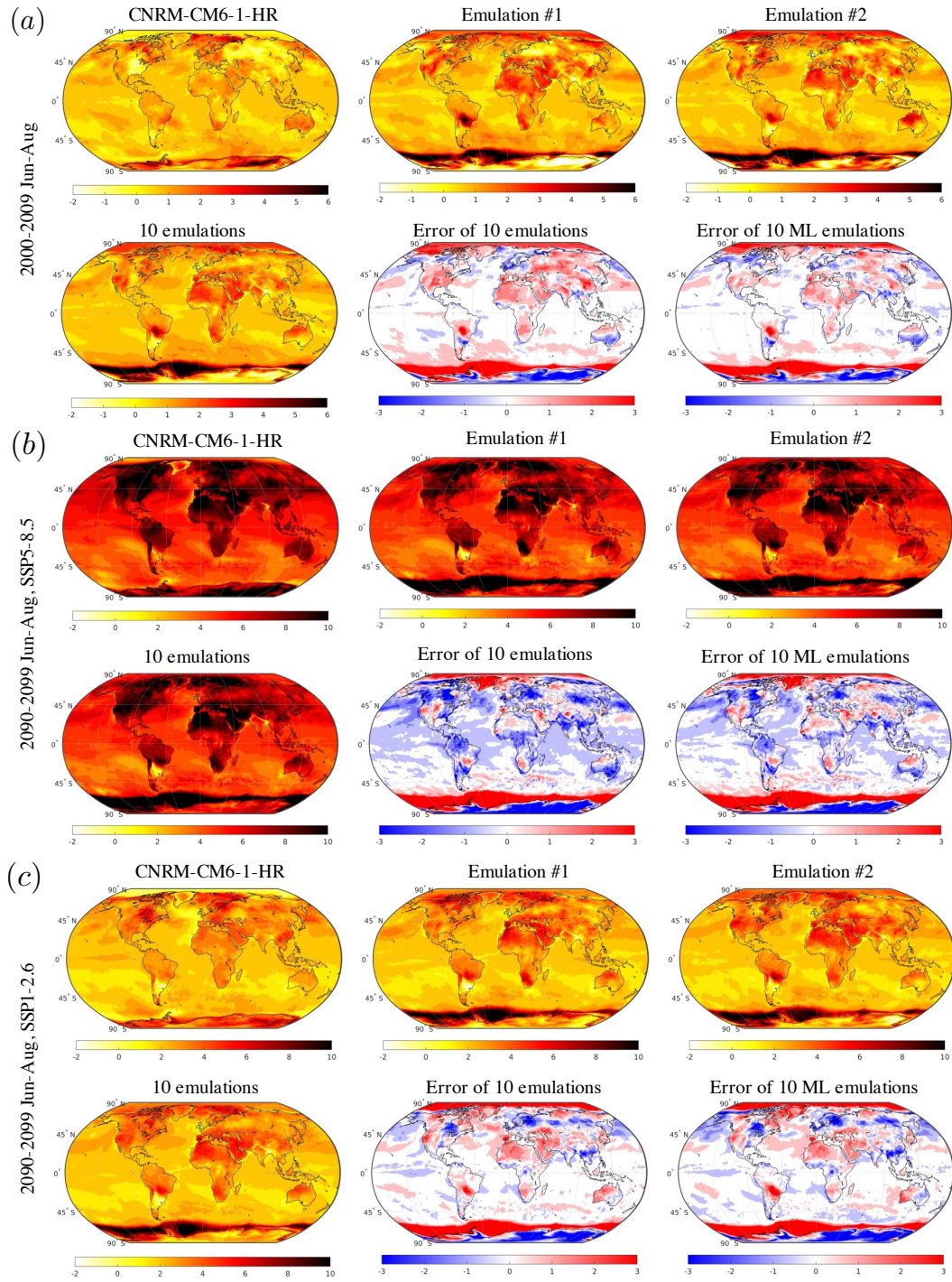
363 is only partially reproduced. These limitations can be alleviated by relaxing the assump-  
 364 tion of the emulator that cross-EOF correlations  $\hat{\mathbf{L}}_s$  are constant, which is explored in  
 365 §4.2. Nonetheless, the underlying climate dynamics, such as the removal of polar am-  
 366 plification due to the loss of ice coverage, is non-linear and non-local, requiring more ju-  
 367 dicious treatment in the construction of emulators. In the SSP1-2.6 scenario (figure 7c),  
 368 changes of standard deviation progress more slowly, and the corresponding emulation  
 369 errors are less severe than in the SSP5-8.5 scenario.

370 We visualize in figure 8 the 97.5% quantile as an example of extreme temperature.  
 371 It is important to note that the baseline temperature for anomalies in figure 8 differs from  
 372 that in figure 6; here, it is based on the 1850-1900 97.5% quantile rather than the 1850-  
 373 1900 average. Within 2000-2009, the emulated quantile (figure 8a) is less accurate than  
 374 the mean (c.f. figure 6a), which is anticipated due to the compounded error from the em-  
 375 ulated standard deviation affecting the quantile estimation. Moreover, the predicted quan-  
 376 tile exhibits greater uncertainty across different emulations, further contaminating the  
 377 accuracy of averaged emulations. In SSP5-8.5 2090-2099 (figure 8b), the increase of quan-  
 378 tile is similar to the mean (figure 6b) at most locations. An interesting trend can be ob-  
 379 served in South Asia: the quantile grows more significantly than the mean in India but  
 380 slightly decreases in Ganges Delta. Since the standard deviation in South Asia remains  
 381 approximately unaffected by the global warming, the change of extreme temperature pre-  
 382 dominantly indicates heavier or thinner tails of the probability distribution. These trends  
 383 are successfully identified by the emulator. The highest error of the emulated quantile  
 384 occurs in Greenland and the Southern Ocean, due to the overestimated standard devi-  
 385 ation as discussed in figure 7. Other error patterns primarily originate from the inter-  
 386 nal variability, as explored by analyzing the temporal evolution of the emulated quan-  
 387 tile from 2010 to SSP5-8.5 2100 (Appendix Appendix A). When applied to the testing  
 388 data under the SSP1-2.6 scenario (figure 8c), the emulator effectively captures the warm-  
 389 ing patterns of extreme temperatures with accuracy comparable to the training data in  
 390 figures 8(a,b). Using the ML model for long-term trends does not improve the quantiles  
 391 of TMX in SSP1-2.6 scenario.

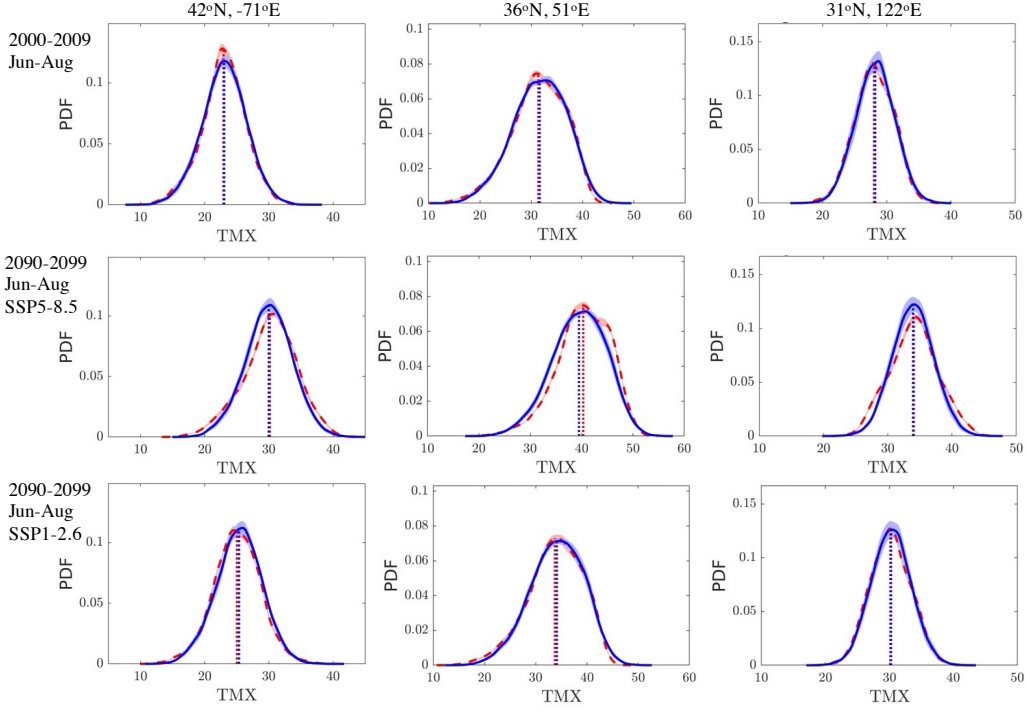
392 The probability density functions of local TMX are plotted in figure 9 at three  $1^\circ \times$   
 393  $1^\circ$  small regions that include major cities: Boston, situated in proximity to the Atlantic  
 394 Ocean; Tehran, featured by the semi-arid climate with hot dry summers; Shanghai, char-  
 395 acterized by the subtropical maritime monsoon climate. All these locations exhibit a sig-  
 396 nificant increase of the extreme temperature in SSP5-8.5 scenario (c.f. figure 8). Over-  
 397 all the emulated PDFs closely match their true profiles, although the deviations in the  
 398 SSP5-8.5 scenario are more appreciable. Since the size of samples (3,680) to estimate the  
 399 true PDF might be insufficient, we quantify the uncertainty by bootstrap resampling,  
 400 as marked by red shaded regions in figure 9. The uncertainty of emulated PDFs are quan-  
 401 tified using one standard deviation of ten emulations, as shown by blue shaded areas.  
 402 Taking the uncertainty of PDFs into consideration, the mismatch between emulated and  
 403 true profiles are less severe. Note that the non-Gaussian shape of the PDF at Tehran  
 404 (middle row in figure 9) is accurately replicated by the emulator, due to the effect of mix-  
 405 ing instantaneous Gaussian TMX with different mean and variance, as discussed at the  
 406 end of §3.2. The accurate emulation of the PDFs demonstrate the capacity of the em-  
 407 ulator to predict any statistics of theoretical and practical interest, including skewness,  
 408 kurtosis, and climate extreme indices.

409 The performance of the emulator in different seasons is examined by the root-mean-  
 410 square error (RMSE) of the statistics and summarized in figure 10. Given a statistic of  
 411 the reference daily maximum temperature  $Q$  and its estimation  $\hat{Q}$ , the RMSE is defined  
 412 as,

$$\text{RMSE} = \left( \frac{1}{S} \int_S (\hat{Q} - Q)^2 \cos \theta d\theta d\varphi \right)^{1/2}. \quad (12)$$

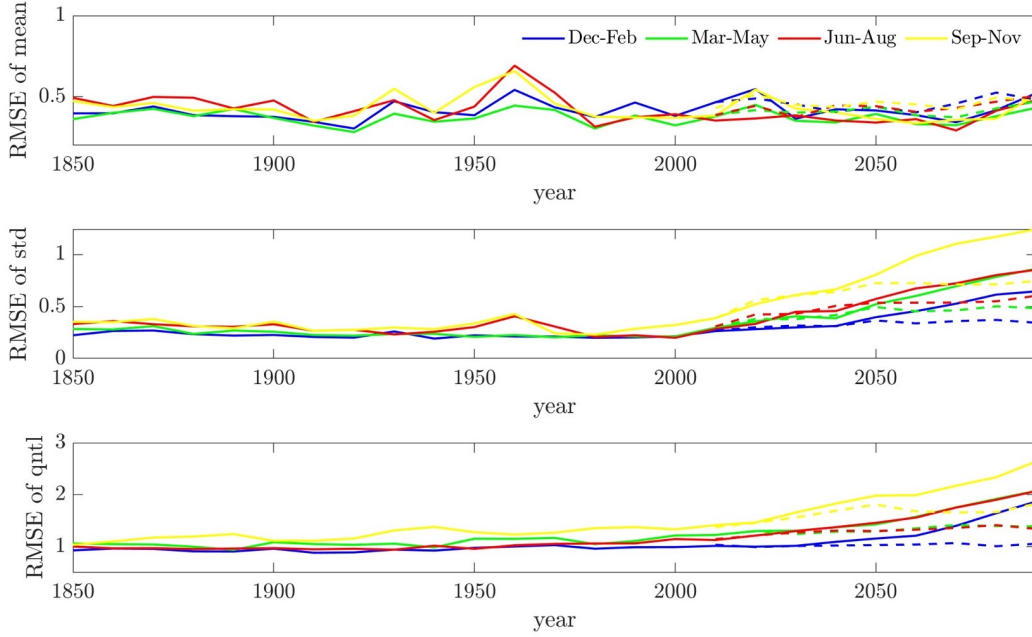


**Figure 8.** Extreme anomaly of Jun-Aug daily maximum temperature, quantified by the 97.5% quantile of local TMX distribution. The quantiles are evaluated using data from (a) 2000-2009, (b) 2090-2099 of the SSP5-8.5 scenario, and (c) 2090-2099 of the SSP1-2.6 scenario. Each sub-figure shows the true mean from CNRM-CM6-1-HR ESM, two sample emulations, average of 10 emulations, error of 10 emulations, and the error of 10 ML emulations. Reference: 1850-1900 Jun-Aug 97.5% quantile of TMX.



**Figure 9.** Probability density function (PDF) of local daily maximum temperature, averaged over three  $1^\circ \times 1^\circ$  regions that include major cities. Left to right columns: Boston ( $42^\circ N, -71^\circ E$ ), Tehran ( $36^\circ N, 51^\circ E$ ) and Shanghai ( $31^\circ N, 122^\circ E$ ). Red dashed line: CNRM-CM6-1-HR simulation data; red shaded region: uncertainty of the true PDF computed by bootstrapping; solid line: 10 emulations; blue shaded region: uncertainty of PDF quantified by one standard deviation of 10 emulations. The PDF are evaluated in decadal windows: (top row) historical, 2000-2009; (middle row) 2090-2099, SSP5-8.5 scenario; (bottom row) 2090-2099, SSP1-2.6 scenario. TMX are shown using degree Celsius.

413 The error in mean TMX remains relatively consistent across seasons and future scenarios.  
 414 Similarly, the standard deviation error is nearly stationary and independent of seasons  
 415 over historical periods. However, in SSP5-8.5 future scenario, seasonal variation becomes  
 416 more pronounced, with the error in Sep-Nov at the end of the century almost doubling  
 417 that of Dec-Feb. The end period of SSP5-8.5 scenario is the most difficult to predict,  
 418 because of the reduced representation accuracy of leading EOFs trained from historical  
 419 data. Additionally, the availability of only a single realization limits the emulator's  
 420 ability to accurately estimate the most extreme warming conditions. The more pronounced  
 421 error in Sep-Nov is due to the more significant influence of global warming on Sep-Nov  
 422 statistics of TMX. Specifically, the Sep-Nov standard deviation of TMX is decreasing  
 423 not only in the Southern Ocean, but also in the Arctic, which are not accurately  
 424 captured by the emulator (see Appendix B for global distribution of standard deviations).  
 425 The SSP1-2.6 future scenario exhibits similar seasonal error variations, albeit with generally  
 426 lower magnitudes compared to SSP5-8.5. Regarding the 97.5% quantiles, their RMSE  
 427 patterns align closely with those observed for the standard deviation, reflecting the same  
 428 underlying climate dynamics. Despite these seasonal variations, the overall error magnitude  
 429 remains relatively consistent across all four seasons throughout the emulated time  
 430 and scenarios, which justifies the application of the emulator across the entire annual cycle.  
 431



**Figure 10.** Root-mean-square error of the mean, standard deviation, and 97.5% quantile of TMX in different seasons. Solid lines: historical and SSP5-8.5 future scenario; dashed lines: SSP1-2.6 future scenario. Blue, green, red, yellow: errors averaged in Dec-Feb, Mar-May, Jun-Aug, Sep-Nov.

432

#### 4.2 Emulation of MPI-ESM1-2-LR large-ensemble dataset

433

434

435

436

When a large ensemble of realizations are available, the assumption of constant cross-mode covariance in the emulator (equation 9) can be relaxed. Specifically, we generalize the emulator of EOF time series (equation 6) by modeling  $\hat{l}_{s,ij}$  as a function of the global mean temperature,

$$\hat{a}_{s,i}(t, \hat{\omega}) = \hat{\mu}_{s,i}(T_{s,g}) + \sum_{j=1}^I \hat{l}_{s,ij}(T_{s,g}) \hat{\eta}_{s,j}(t, \hat{\omega}), \quad i = 1, 2, \dots, I. \quad (13)$$

437

438

439

440

In order to estimate the relation between  $\hat{l}_{s,ij}$  and  $T_{s,g}$ , we follow similar procedures as §3.2. Given the true EOF time series  $\mathbf{a}(t)$ , we remove the linear trends of seasonal mean  $\hat{\mu}_s(T_{s,g})$ , compute the covariance of  $\mathbf{a}_s - \hat{\mu}_s$  in each year, and perform Cholesky decomposition of the covariance matrix,

$$\bar{\Sigma}_s(t) = \left\langle (\mathbf{a}_s - \hat{\mu}_s) (\mathbf{a}_s - \hat{\mu}_s)^\top \right\rangle_{s\omega}, \quad \bar{\Sigma}_s(t) = \bar{\mathbf{L}}_s(t) \bar{\mathbf{L}}_s^\top(t), \quad (14)$$

441

442

443

444

445

where  $\langle \cdot \rangle_{s\omega}$  denotes an average over the ensemble and season  $s$  in each year. An intuitive but risky idea is modelling each entry of  $\bar{\Sigma}_s(t)$  as a linear function of the global mean temperature. Such a strategy cannot guarantee the positive definite property of the estimated covariance matrix. This limitation can be overcome by modelling  $\bar{\mathbf{L}}_s(t)$  as linear functions of  $T_{s,g}$ ,

$$\hat{\mathbf{L}}_s(T_{s,g}) = \hat{\mathbf{P}}_{s,0} + T_{s,g} \hat{\mathbf{P}}_{s,1}. \quad (15)$$

446

447

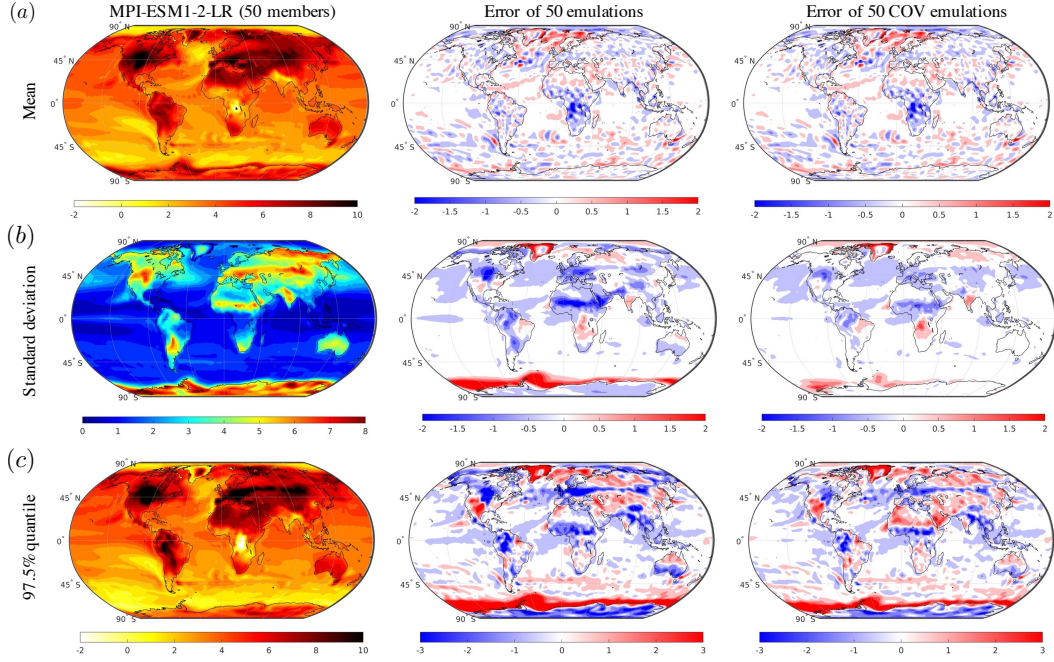
448

449

Since  $\hat{\mathbf{L}}_s(T_{s,g})$  is lower triangular,  $\hat{\mathbf{P}}_{s,0}$  and  $\hat{\mathbf{P}}_{s,1}$  inherit this property, and each of their non-zero entries is computed by the method of least squares. Multiplying  $\mathbf{a}_s - \hat{\mu}_s$  by  $\hat{\mathbf{L}}_s^{-1}(T_{s,g})$ , we can extract the time series that are approximately uncorrelated in each season of each year,

$$\boldsymbol{\eta}_s(t, \omega) = \hat{\mathbf{L}}_s^{-1}(T_{s,g}) (\mathbf{a}_s(t, \omega) - \hat{\mu}_s). \quad (16)$$





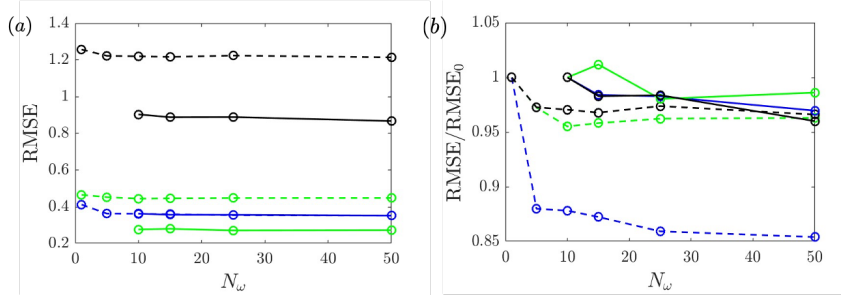
**Figure 11.** Statistics of Jun-Aug daily maximum temperature of MPI-ESM1-2-LR dataset and the emulations. All the statistics are evaluated in 2090-2099 of the SSP5-8.5 scenario. (a) Mean anomaly from 1850-1900; (b) Standard deviation; (c) Anomaly of 97.5% quantile of local TMX distribution from the 1850-1900 value.

450 The autocorrelation of each component of  $\eta_s$  will be used to generate Gaussian processes.  
 451 The remaining procedures for constructing the emulator are the same as in §3.2 and there-  
 452 fore not repeated here for conciseness.

453 Although the generalization introduced in (13-15) has the potential to improve the  
 454 performance of the emulator, it is only applicable when the data are sufficient to obtain  
 455 converged time-dependent covariance matrices. A minimum requirement for the amount  
 456 of data is that the number of samples for computing the covariance matrix (14) must  
 457 exceed the number of EOFs, or equivalently the size of  $\bar{\Sigma}_s(t)$ . This requirement is not  
 458 satisfied by the CNRM-CM6-1-HR dataset. For example, in Northern Hemisphere sum-  
 459 mer of every year we have 92 samples to compute  $\bar{\Sigma}_s(t)$ , but the number of EOFs used  
 460 in the emulator is 2,000. As a result, the computed covariance matrix is not even full  
 461 rank, consisting of spurious correlations that contaminate the dependence on time or global  
 462 mean temperature.

463 To distinguish from the emulator introduced in §3.2, all the results generated us-  
 464 ing (13-15) will be termed as COV emulations. Both types of emulators are applied to  
 465 the MPI-ESM1-2-LR dataset to compare their performance. Different from the CNRM-  
 466 CM6-1-HR dataset that requires 2,000 EOFs to represent 95% of the total variance, only  
 467 1,000 EOFs are sufficient to model the MPI-ESM1-2-LR dataset due to lower spatial res-  
 468 olutions. All the 50 realizations of the historical and SSP5-8.5 scenarios are used to com-  
 469 pute the EOFs and train the stochastic emulators of the EOF time series.

470 Since the error of emulated statistics were highest in SSP585 2090-2099 for the CNRM-  
 471 CM6-1-HR dataset, we focus on this time window to compare the performance of the  
 472 emulators. The results are visualized in figure 11. Overall the warming trend predicted  
 473 by MPI-ESM1-2-LR model is less pronounced than the CNRM-CM6-1-HR model, which



**Figure 12.** (a) Absolute and (b) Relative root-mean-square error of emulated statistics versus the number of realizations used for training the emulator. The statistics are evaluated in SSP5-8.5 2090-2099 Jun-Aug. Dashed lines: error of 50 emulations; Solid lines: error of 50 COV emulations. Blue, green, black: error of the mean, standard deviation, and 97.5% quantile. The relative errors in (b) are normalized by the values associated with the smallest  $N_\omega$ .

474 is consistent with previous studies on equilibrium climate sensitivity of ESMs (Tokarska  
 475 et al., 2020). In figure 11a, the error of the mean anomaly of both emulators are almost  
 476 identical, which is expected since the same linear model is adopted for the seasonal mean  
 477 of EOF coefficients. The error of local standard deviation, as shown in panel b, is sig-  
 478 nificantly reduced by modeling the variations of covariance matrix. For example, the high-  
 479 est errors in North Africa and the Southern Ocean are decreased by approximately  $2^\circ C$ ,  
 480 which confirms the speculation in §4.1 that these errors are mostly associated with time-  
 481 dependent cross-mode correlations. As a result of more accurate estimation of local vari-  
 482 ance in COV emulations, the quantiles in panel c are also reproduced with lower errors.

483 To assess the influence of ensemble size of the training data on both emulators, we  
 484 calculated the root-mean-square error (RMSE) of the emulated statistics. The results  
 485 are reported in figure 12 for the mean, standard deviation, and 97.5% quantile in SSP5-  
 486 8.5 2090-2099 Jun-Aug. When  $N_\omega$  realizations are available for training the emulator,  
 487 the true statistics  $\mathcal{Q}$  are also evaluated using the same  $N_\omega$  realizations, while the em-  
 488 ulators are always performed 50 times to generate converged statistics,  $\hat{\mathcal{Q}}$ . In panel a,  
 489 compared with the constant-covariance emulator (dashed lines), the COV emulator (solid  
 490 lines) achieves approximately 40% error reduction in the standard deviation and 30%  
 491 in the quantile. However, the COV emulator requires at least ten realizations to ensure  
 492 the positive definiteness of the covariance matrices. To highlight the dependence of em-  
 493 ulation error on the ensemble size  $N_\omega$ , the RMSE is normalized by the value associated  
 494 with the smallest  $N_\omega$  attempted. The results are shown in figure 12b. For the constant-  
 495 covariance emulator (dashed lines), as the size of ensemble is increased from one to ten,  
 496 the RMSE of mean, standard deviation and quantile are respectively decreased by 12%,  
 497 4.5% and 3.0%. These error reductions suggest that the emulation accuracy is generally  
 498 improved when the impact of climate internal variability is alleviated in the training data.  
 499 Such a trend is also consistent with conclusions of previous studies (Tebaldi et al., 2021)  
 500 that approximately ten realizations are required to capture the ensemble variance ac-  
 501 curately. As the ensemble size reaches 50, further error reduction becomes negligible for  
 502 the standard deviation (green dashed) and quantile (black dashed), suggesting dimin-  
 503 ishing returns from larger training datasets. In contrast, the COV emulator shows con-  
 504 tinued improvement, with a reduction in error of 1.4% for the standard deviation and  
 505 4.0% for the quantile, since larger-ensemble data can still help improve the emulated co-  
 506 variance matrices. Despite these gains, the COV emulator constructed with ten ensemble  
 507 members already provides an accurate estimation of the statistics of extreme tem-  
 508 perature. These results indicate that as long as the amount of training data are suffi-

509 cient to construct the COV emulator, the performance of the emulator is robust against  
 510 the ensemble size of realizations.

## 511 5 Conclusions and Discussion

512 We have developed a framework of a spatially resolved stochastic emulator that es-  
 513 timates the full statistics of climate extremes. The emulator was trained and tested us-  
 514 ing the daily maximum temperature data from CNRM-CM6-1-HR and MPI-ESM1-2-  
 515 LR Earth system simulations in CMIP6. To reduce the dimensionality of the global cli-  
 516 mate system and achieve speedy emulations, we extract empirical orthogonal functions  
 517 of daily maximum temperature data and assume their shapes remain unchanged across  
 518 different climate change scenarios. The time series of EOF coefficients are decomposed  
 519 as the combination of long term trends of seasonal statistics and conditionally Gaussian  
 520 daily fluctuations. The former, including seasonal mean and variance, are approximated  
 521 as linear or machine-learned functions of the global mean temperature, while the daily  
 522 fluctuations are modeled as Gaussian autoregressive processes that are scaled by the cross  
 523 correlations of different EOFs. While the statistics of the emulator, conditioned on sea-  
 524 son and global mean temperature, are assumed to be Gaussian, the long term statistics  
 525 of the model do not produce normal distribution due to variation of the global mean tem-  
 526 perature. However, the possibility of heavy tailed daily temperature fluctuations is not  
 527 covered and is left for future work.

528 The performance of the emulator is evaluated on the CNRM-CM6-1-HR dataset  
 529 due to its high spatial resolution. Trained on historical and SSP5-8.5 scenario, the em-  
 530 ulated time series accurately reproduce the evolution of the seasonal mean and the Fourier  
 531 spectra of daily fluctuations. After generating the spatiotemporal evolution of the in-  
 532 stantaneous daily maximum temperature, the emulator’s performance is systematically  
 533 evaluated on the ten-year Jun-Aug statistics, including the mean, standard deviation,  
 534 quantile, and the full probability density function. Remarkably, the emulator reproduces  
 535 the quantile anomaly in response to climate change and effectively captures the non-Gaussian  
 536 profiles of the local PDF. When tested on the SSP1-2.6 scenario that is not included in  
 537 the training data, the full statistics are also accurately predicted, which demonstrates  
 538 the potential of the emulator to be applied to various climate change scenarios. While  
 539 using neural networks to represent the impact of global warming improves the emula-  
 540 tor’s performance on the training SSP5-8.5 scenario compared to linear functions, this  
 541 improvement does not extend to the SSP1-2.6 scenario used for validation.

542 Based on MPI-ESM1-2-LR large-ensemble datasets, we further developed the em-  
 543 ulator by modelling the variation of the cross-mode covariance as linear functions of the  
 544 global mean temperature. Such a refinement helps reduce the root-mean-square error  
 545 of emulated local statistics by 50%. By progressively increasing the number of ensem-  
 546 ble members in the training data, we assessed the impact of climate internal variabil-  
 547 ity on performance of both emulators. Overall the RMSE of statistics decrease with larger  
 548 ensemble. When more than ten members are included, the accuracy of the constant-covariance  
 549 emulator approximately saturates, but COV emulator shows continued improvement. As  
 550 long as there are sufficient training data to construct the COV emulator, its performance  
 551 remains relatively stable regardless of the ensemble size of realizations.

552 There are numerous pathways for generalizing the emulator to further improve its  
 553 accuracy, and we outline a few possibilities below. First, the time-lagged covariance be-  
 554 tween different EOFs can be included into the emulator to achieve a better estimation  
 555 of the full probability distribution of local temperature (Wan et al., 2021). Second, in-  
 556 stead of using the global mean temperature as the driver, the emulator can be param-  
 557 eterized using the emission history of greenhouse gases, the equivalent radiative forcing,  
 558 or aerosol concentrations (Castruccio et al., 2014; Freese et al., 2024). Such an exten-  
 559 sion will take into account the memory effect and facilitate the application of the em-

560 ulator into scenarios where the evolution of global mean temperature is non-monotonic.  
 561 Third, the Empirical Orthogonal Functions can be replaced by more state-of-the-art deep  
 562 learning methods, such as Autoencoders, to nonlinearly reduce the dimensionality of the  
 563 climate system (Kramer, 1991). Lastly, a recently proposed non-intrusive machine-learning  
 564 framework shows promise for further improving the emulator’s accuracy (Barthel Sorensen  
 565 et al., 2024). This approach focuses on learning a debiasing operator that takes the em-  
 566 ulated time series of temperature fields as input and corrects them to better match the  
 567 reference data from ESMs. Once trained on a few scenarios, this debiasing operator can  
 568 be applied to correct the emulations in other unseen climate change scenarios. Despite  
 569 these potential enhancements, the emulator successful estimation of extreme tempera-  
 570 ture statistics is promising and suggests its applicability to other variables, such as hu-  
 571 midity, precipitation, and wind speed, which will better assist with risk management of  
 572 climate extremes.

## 573 **Appendix A Temporal evolution of emulated quantile in SSP5-8.5 sce-** 574 **nario**

575 In this appendix, we provide more details about the temporal evolution of statis-  
 576 tics of extreme temperature in SSP5-8.5 scenario. Similar to figure 8, we evaluate the  
 577 97.5% quantile of the local TMX using ten-year Jun-Aug data. The anomaly of quan-  
 578 tiles against 1850-1900 reference are visualized in figure A1 and A2 from 2010 to 2089.  
 579 Overall the regions with the most rapid increase of extreme temperature are correctly  
 580 identified by the emulator. Two categories of error patterns can be observed. The first  
 581 type is relatively independent of time, such as the overestimated quantile in Greenland.  
 582 The second type is more stochastic, sometime even changing signs across different time  
 583 windows, such as the North America and southern Africa. These error patterns are prob-  
 584 ably associated with the internal variability of the global climate system and require more  
 585 realizations of the Earth system simulations to converge.

## 586 **Appendix B Emulated statistics in other seasons**

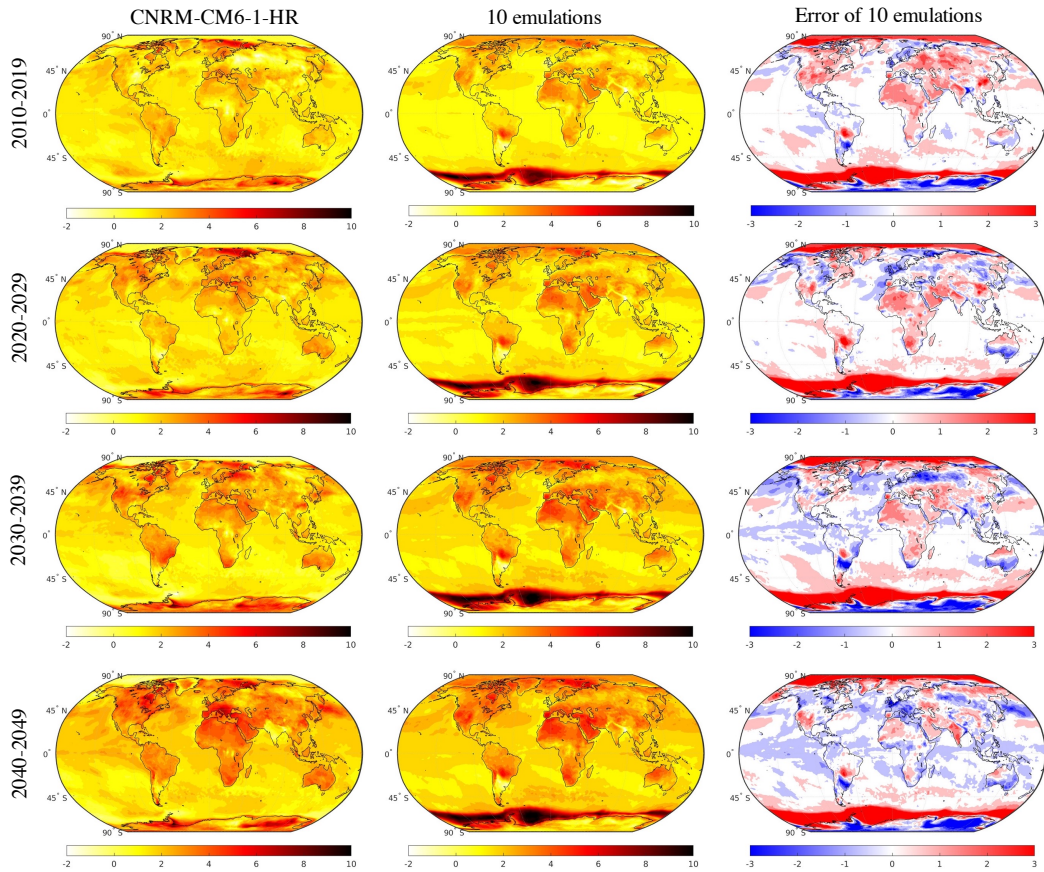
587 This appendix presents the statistics of TMX across different seasons and their cor-  
 588 responding emulation errors. The local standard deviation in 2090-2099 of the SSP5-8.5  
 589 scenario is shown in figure B1. In Dec-Feb, the error reaches its maximum in the Arc-  
 590 tic, contrasting with the Jun-Aug pattern where the error peaks in the Southern Ocean  
 591 (c.f. figure B1). This seasonal difference is likely associated with the sea ice coverage.  
 592 During Dec-Feb, Antarctic sea ice consistently retreats almost to the coastline in both  
 593 historical and global warming scenarios. Therefore, the standard deviation of TMX in  
 594 this season is less affected by warming conditions compared to Jun-Aug. Mar-May and  
 595 Sep-Nov present a more complex picture. During these transitional seasons, sea ice cov-  
 596 erage in both polar regions is highly sensitive to climate change. The emulator strug-  
 597 gles to capture the associated trends in standard deviations, resulting in high errors in  
 598 these areas. The error patterns of 97.5% quantile are analogous to the standard devi-  
 599 ation, as shown in figure B2.

## 600 **Open Research**

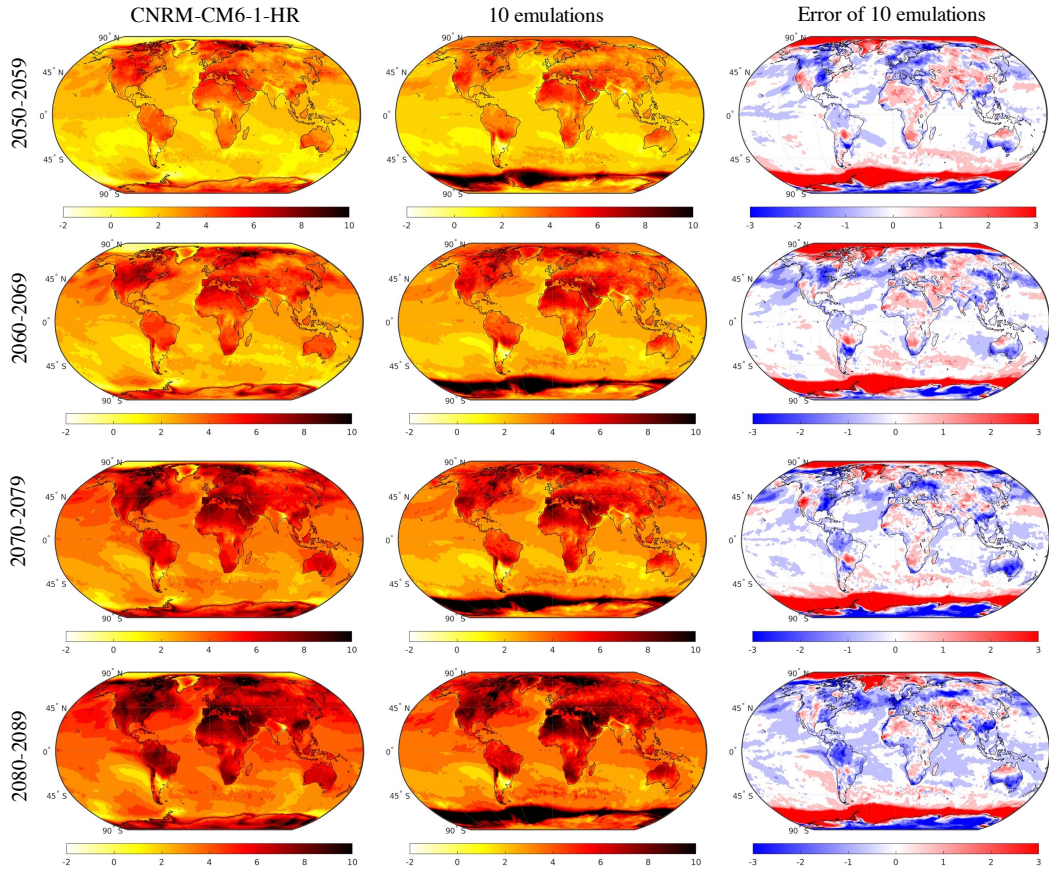
601 All code to reproduce this work is available at [https://github.com/mzwang2012/](https://github.com/mzwang2012/SEM.TMX.git)  
 602 [SEM.TMX.git](https://github.com/mzwang2012/SEM.TMX.git). The raw data from CMIP6 were retrieved through the Earth System Grid  
 603 Federation interface <https://aims2.llnl.gov/search/cmip6/>.

## 604 **Acknowledgments**

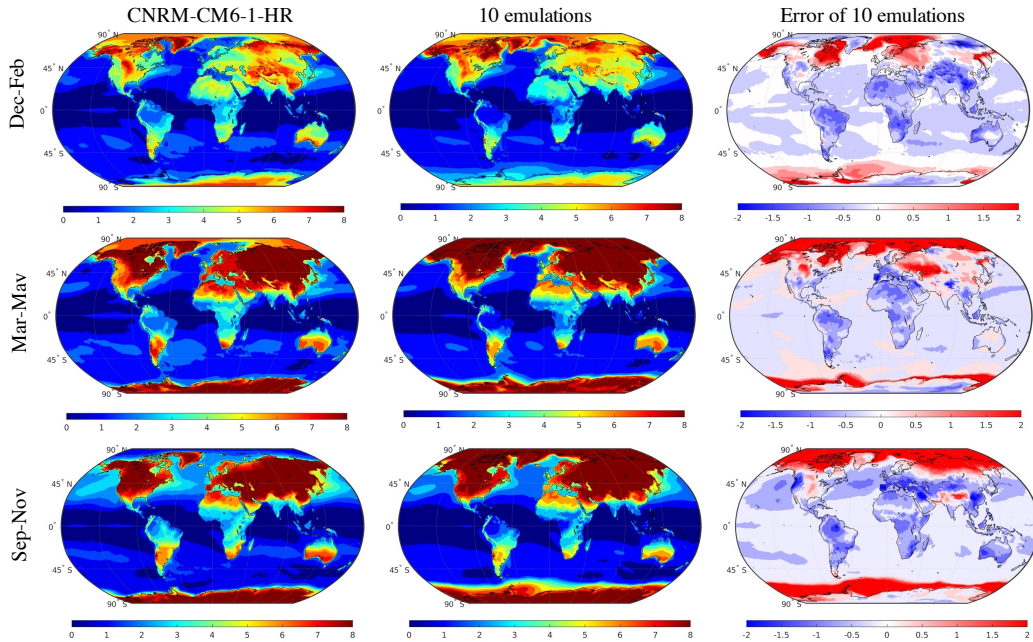
605 This research was part of the Bringing Computation to the Climate Challenge (BC3)  
 606 project and supported by Schmidt Sciences through the MIT Climate Grand Challenges.



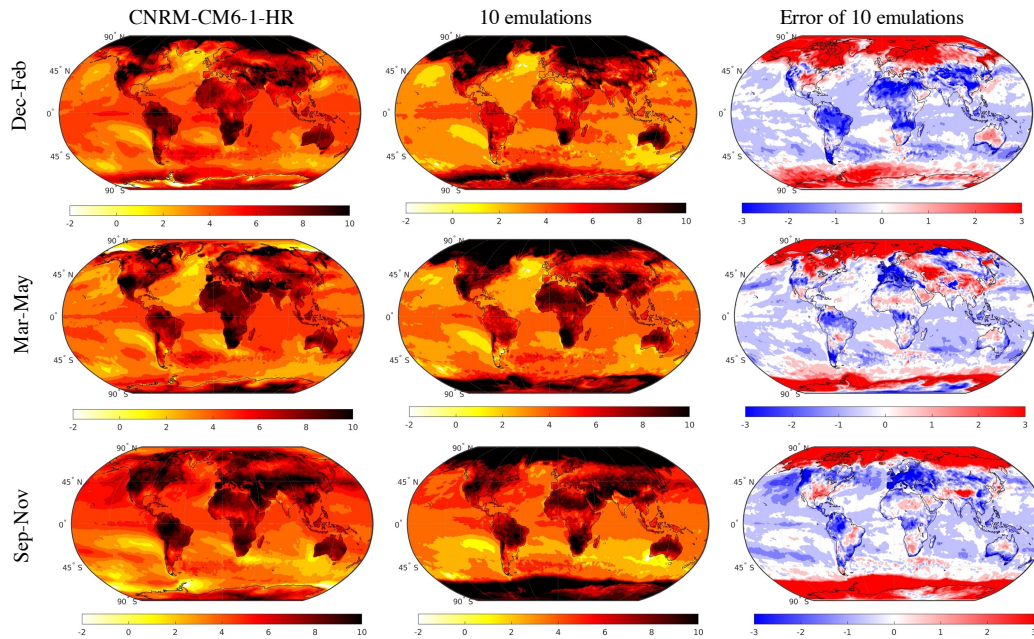
**Figure A1.** Extreme anomaly of ten-year Jun-Aug daily maximum temperature, quantified by the 97.5% quantile of local TMX distribution. The quantiles are evaluated for SSP5-8.5 scenario within 2010-2019, 2020-2029, 2030-2039, 2040-2049, respectively. Reference: 1850-1900 Jun-Aug 97.5% quantile of TMX.



**Figure A2.** Same as figure A1, but shown for 2050-2059, 2060-2069, 2070-2079, 2080-2089, respectively.



**Figure B1.** Standard deviation of ten-year seasonal daily maximum temperature, evaluated for Dec-Feb, Mar-May, and Sep-Nov in 2090-2099 of the SSP5-8.5 future scenario.



**Figure B2.** Extreme anomaly of ten-year seasonal daily maximum temperature, quantified by the 97.5% quantile of local TMX distribution. The quantiles are evaluated for Dec-Feb, Mar-May, and Sep-Nov in 2090-2099 of the SSP5-8.5 future scenario.. Reference: 1850-1900 97.5% quantile of TMX of each season.

## References

607

608 Alexeeff, S. E., Nychka, D., Sain, S. R., & Tebaldi, C. (2018). Emulating mean  
609 patterns and variability of temperature across and within scenarios in anthro-  
610 pogenic climate change experiments. *Climatic Change*, *146*, 319–333.

611 Allen, S., Barros, V., (Canada, I., (UK, D., Cardona, O., Cutter, S., ... (USA, T.  
612 (2012, nov). *Managing the Risks of Extreme Events and Disasters to Advance  
613 Climate Change Adaptation. Special Report of Working Groups I and II of the  
614 Intergovernmental Panel on Climate Change.* doi: 10.13140/2.1.3117.9529

615 Amaya, D. J. (2019). The pacific meridional mode and enso: A review. *Current Cli-  
616 mate Change Reports*, *5*(4), 296–307.

617 AON. (2020). *Weather, climate & catastrophe insight, 2020 annual report.*

618 Arbabi, H., & Sapsis, T. (2022). Generative stochastic modeling of strongly non-  
619 linear flows with non-gaussian statistics. *SIAM/ASA Journal on Uncertainty  
620 Quantification*, *10*(2), 555–583.

621 Barriopedro, D., Fischer, E. M., Luterbacher, J., Trigo, R. M., & García-Herrera, R.  
622 (2011). The hot summer of 2010: Redrawing the temperature record map of  
623 europe. *Science*, *332*(6026), 220-224.

624 Barthel Sorensen, B., Charalampopoulos, A., Zhang, S., Harrop, B., Leung, L., &  
625 Sapsis, T. P. (2024). A non-intrusive machine learning framework for de-  
626 biasing long-time coarse resolution climate simulations and quantifying rare  
627 events statistics. *Journal of Advances in Modeling Earth Systems*, *16*(3),  
628 e2023MS004122.

629 Beusch, L., Gudmundsson, L., & Seneviratne, S. I. (2020). Emulating earth system  
630 model temperatures with mesmer: from global mean temperature trajecto-  
631 ries to grid-point-level realizations on land. *Earth System Dynamics*, *11*(1),  
632 139–159.

- 633 Castruccio, S., McInerney, D. J., Stein, M. L., Crouch, F. L., Jacob, R. L., & Moyer,  
634 E. J. (2014). Statistical emulation of climate model projections based on  
635 precomputed gcm runs. *Journal of Climate*, *27*(5), 1829–1844.
- 636 Fogt, R. L., & Marshall, G. J. (2020). The southern annular mode: variability,  
637 trends, and climate impacts across the southern hemisphere. *Wiley Interdisci-*  
638 *plinary Reviews: Climate Change*, *11*(4), e652.
- 639 Freese, L. M., Fiore, A. M., & Selin, N. E. (2024). Spatially resolved temperature re-  
640 sponse functions to co2 emissions. *Authorea Preprints*.
- 641 Gao, Y., Leung, L. R., Lu, J., & Masato, G. (2015). Persistent cold air outbreaks  
642 over north america in a warming climate. *Environmental Research Letters*,  
643 *10*(4), 044001.
- 644 Geogdzhayev, G., Souza, A., Ferrari, R., & Flierl, G. R. (2024). *A statistical emula-*  
645 *tor design for averaged climate fields*. (Personal Communications)
- 646 Hannachi, A., Jolliffe, I. T., & Stephenson, D. B. (2007). Empirical orthogonal func-  
647 tions and related techniques in atmospheric science: A review. *International*  
648 *Journal of Climatology: A Journal of the Royal Meteorological Society*, *27*(9),  
649 1119–1152.
- 650 Huntingford, C., Jones, P. D., Livina, V. N., Lenton, T. M., & Cox, P. M. (2013).  
651 No increase in global temperature variability despite changing regional pat-  
652 terns. *Nature*, *500*(7462), 327–330.
- 653 Kaltenborn, J., Lange, C., Ramesh, V., Brouillard, P., Gurwicz, Y., Nagda, C., ...  
654 Rolnick, D. (2023). Climateset: A large-scale climate model dataset for ma-  
655 chine learning. In A. Oh, T. Neumann, A. Globerson, K. Saenko, M. Hardt, &  
656 S. Levine (Eds.), *Advances in neural information processing systems* (Vol. 36,  
657 pp. 21757–21792). Curran Associates, Inc.
- 658 Kalvová, J., & Nemesšová, I. (1998). Estimating autocorrelations of daily extreme  
659 temperatures in observed and simulated climates. *Theoretical and applied cli-*  
660 *matology*, *59*, 151–164.
- 661 Kapmeier, F., Greenspan, A., Jones, A., & Sterman, J. (2021). Science-based  
662 analysis for climate action: how hsbc bank uses the en-roads climate policy  
663 simulation. *System dynamics review: the journal of the System Dynamics*  
664 *Society*, *37*(4), 333–352.
- 665 Kramer, M. A. (1991). Nonlinear principal component analysis using autoassociative  
666 neural networks. *AICHE journal*, *37*(2), 233–243.
- 667 Link, R., Snyder, A., Lynch, C., Hartin, C., Kravitz, B., & Bond-Lamberty, B.  
668 (2019). Fldgen v1. 0: an emulator with internal variability and space–time  
669 correlation for earth system models. *Geoscientific Model Development*, *12*(4),  
670 1477–1489.
- 671 Lorenz, E. N. (1956). *Empirical orthogonal functions and statistical weather predic-*  
672 *tion* (Vol. 1). Massachusetts Institute of Technology, Department of Meteorol-  
673 ogy Cambridge.
- 674 Lütjens, B., Ferrari, R., Watson-Parris, D., & Selin, N. (2024). The impact of in-  
675 ternal variability on benchmarking deep learning climate emulators. *arXiv*  
676 *preprint arXiv:2408.05288*.
- 677 Meehl, G. A., & Tebaldi, C. (2004). More intense, more frequent, and longer lasting  
678 heat waves in the 21st century. *Science*, *305*(5686), 994–997.
- 679 Meinshausen, M., Raper, S. C., & Wigley, T. M. (2011). Emulating coupled  
680 atmosphere-ocean and carbon cycle models with a simpler model, magicc6–  
681 part 1: Model description and calibration. *Atmospheric Chemistry and*  
682 *Physics*, *11*(4), 1417–1456.
- 683 Mitchell, T. D. (2003). Pattern scaling: an examination of the accuracy of the tech-  
684 nique for describing future climates. *Climatic change*, *60*(3), 217–242.
- 685 Mohamad, M. A., & Sapsis, T. P. (2015). Probabilistic description of extreme events  
686 in intermittently unstable dynamical systems excited by correlated stochastic  
687 processes. *SIAM/ASA Journal on Uncertainty Quantification*, *3*(1), 709–736.



- 688 Nicholls, Z., Meinshausen, M., Lewis, J., Corradi, M. R., Dorheim, K., Gasser, T.,  
689 ... others (2021). Reduced complexity model intercomparison project phase  
690 2: Synthesizing earth system knowledge for probabilistic climate projections.  
691 *Earth's Future*, 9(6), e2020EF001900.
- 692 Nicholls, Z. R., Meinshausen, M., Lewis, J., Gieseke, R., Dommenges, D., Dorheim,  
693 K., ... others (2020). Reduced complexity model intercomparison project  
694 phase 1: Protocol, results and initial observations. *Geoscientific Model Devel-*  
695 *opments*.
- 696 Olonscheck, D., & Notz, D. (2017). Consistently estimating internal climate variabil-  
697 ity from climate model simulations. *Journal of Climate*, 30(23), 9555–9573.
- 698 Percival, D. B. (1993). Simulating gaussian random processes with specified spectra.  
699 *Computing Science and Statistics*, 534–534.
- 700 Quilcaille, Y., Gudmundsson, L., Beusch, L., Hauser, M., & Seneviratne, S. I.  
701 (2022). Showcasing mesmer-x: Spatially resolved emulation of annual max-  
702 imum temperatures of earth system models. *Geophysical Research Letters*,  
703 49(17), e2022GL099012.
- 704 Räisänen, J. (2002). Co2-induced changes in interannual temperature and precip-  
705 itation variability in 19 cmip2 experiments. *Journal of Climate*, 15(17), 2395–  
706 2411.
- 707 Reed, K. A., Wehner, M. F., & Zarzycki, C. M. (2022). Attribution of 2020 hurri-  
708 cane season extreme rainfall to human-induced climate change. *Nature commu-*  
709 *nications*, 13(1), 1905.
- 710 Rooney-Varga, J. N., Hensel, M., McCarthy, C., McNeal, K., Norfles, N., Rath, K.,  
711 ... Serman, J. D. (2021). Building consensus for ambitious climate action  
712 through the world climate simulation. *Earth's Future*, 9(12), e2021EF002283.
- 713 Seneviratne, S. I., Donat, M. G., Pitman, A. J., Knutti, R., & Wilby, R. L. (2016).  
714 Allowable co2 emissions based on regional and impact-related climate targets.  
715 *Nature*, 529(7587), 477–483.
- 716 Sirovich, L. (1987). Turbulence and the dynamics of coherent structures. i. coherent  
717 structures. *Quarterly of applied mathematics*, 45(3), 561–571.
- 718 Taira, K., Hemati, M. S., Brunton, S. L., Sun, Y., Duraisamy, K., Bagheri, S., ...  
719 Yeh, C.-A. (2020). Modal analysis of fluid flows: Applications and outlook.  
720 *AIAA journal*, 58(3), 998–1022.
- 721 Tebaldi, C., Armbruster, A., Engler, H., & Link, R. (2020). Emulating climate ex-  
722 treme indices. *Environmental Research Letters*, 15(7), 074006.
- 723 Tebaldi, C., Dorheim, K., Wehner, M., & Leung, R. (2021). Extreme metrics from  
724 large ensembles: investigating the effects of ensemble size on their estimates.  
725 *Earth System Dynamics*, 12(4), 1427–1501.
- 726 Thompson, D. W., & Wallace, J. M. (1998). The arctic oscillation signature in the  
727 wintertime geopotential height and temperature fields. *Geophysical research*  
728 *letters*, 25(9), 1297–1300.
- 729 Tokarska, K. B., Stolpe, M. B., Sippel, S., Fischer, E. M., Smith, C. J., Lehner, F.,  
730 & Knutti, R. (2020). Past warming trend constrains future warming in cmip6  
731 models. *Science advances*, 6(12), eaaz9549.
- 732 Wallace, J. M., & Gutzler, D. S. (1981). Teleconnections in the geopotential height  
733 field during the northern hemisphere winter. *Monthly weather review*, 109(4),  
734 784–812.
- 735 Wan, Z. Y., Dodov, B., Lessig, C., Dijkstra, H., & Sapsis, T. P. (2021). A data-  
736 driven framework for the stochastic reconstruction of small-scale features with  
737 application to climate data sets. *Journal of Computational Physics*, 442,  
738 110484.
- 739 Watson-Parris, D., Rao, Y., Olivie, D., Seland, Ø., Nowack, P., Camps-Valls, G.,  
740 ... others (2022). Climatebench v1. 0: A benchmark for data-driven cli-  
741 mate projections. *Journal of Advances in Modeling Earth Systems*, 14(10),  
742 e2021MS002954.

743 Wehner, M., Gleckler, P., & Lee, J. (2020). Characterization of long period return  
744 values of extreme daily temperature and precipitation in the cmip6 models:  
745 Part 1, model evaluation. *Weather and Climate Extremes*, *30*, 100283.