

Lea Beusch
Institute for Atmospheric and
Climate Science ETH Zurich
Universitaetstrasse 16
CH-8092 Zurich
Phone: +41 44 633 36 24
E-mail: lea.beusch@env.ethz.ch



Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

Earth System Dynamics Editorial Board

Zurich, 08 November 2019

Emulating Earth System Model temperatures: from global mean temperature trajectories to grid-point level realizations on land

Dear Dr. Krishnan,

Please find enclosed the revised version of our manuscript and additionally a version with tracked changes directly attached to this letter. For the answer to the reviewers we refer to the point-by-point answers we posted in the interactive discussion.

Based on the reviewers suggestions and your comments, the following main changes were made to the manuscript:

1. A new terminology has been introduced in which climate model runs used to train the emulator are referred to as training runs and the additional initial-condition ensemble members are referred to as test runs. The new terminology improves the clarity of the manuscript and makes it more apparent that we evaluate our emulations both with respect to the training run and independent initial-condition ensemble members which were not employed during training.
2. The quantitative verification section of the manuscript has been extended substantially. In addition to the regional-scale emulation verification, a separate verification of local trends and local variability is now included with a special focus set on the space-time characteristics of the local variability and on how distinguishable climate model runs are from emulations.
3. A supplementary comparison between emulations and simple pattern scaling results has been added to demonstrate their differences.
4. The added value of our study compared to existing literature is now explicitly stated and the extent to which previously developed methods are used in this study is discussed more extensively.

Please also note that during the review process we have improved the local residual temperature variability module in our emulator. The new module locally fits AR(1) processes instead of AR(p) processes which makes more time slots available to estimate the spatial cross-correlations between grid points and therefore decreases the amount of regularization needed.

For all additional changes we implemented based on the reviews, please refer directly to the revised manuscript.

We are confident that the revisions increased the quality of the manuscript.

Yours sincerely,

A handwritten signature in blue ink that reads 'L. Beusch'.

Lea Beusch
(on behalf of all co-authors)

Emulating Earth System Model temperatures: from global mean temperature trajectories to grid-point level realizations on land

Lea Beusch¹, Lukas Gudmundsson¹, and Sonia I. Seneviratne¹

¹Institute for Atmospheric and Climate Science, ETH Zurich, Zurich, Switzerland

Correspondence: Lea Beusch (lea.beusch@env.ethz.ch)

Abstract. Earth System Models (ESMs) are invaluable tools to study the climate system’s response to specific greenhouse gas emission pathways. Large single-model initial-condition and multi-model ensembles are used to investigate the range of possible responses and serve as input to climate impact and integrated assessment models. Thereby, climate signal uncertainty is propagated along the uncertainty chain and its effect on interactions between humans and the Earth system can be quantified. However, generating both single-model initial-condition and multi-model ensembles is computationally expensive. In this study, we assess the feasibility of geographically-explicit climate model emulation, i.e., of statistically producing large ensembles of global ^[.1]land temperature field time series that closely resemble ESM runs at a negligible computational cost^[.2]. For this purpose, we develop a modular ^[.3]emulation framework which consists of (i) a global mean temperature ^[.4]module, (ii) a local ^[.5]temperature response module, and (iii) a local residual temperature variability ^[.6]module. We first show that to successfully mimic single-model initial-condition ensembles of yearly temperature from 1870 to 2100 on grid-point to regional scales, it is sufficient to train on a single ESM run, but separate emulators need to be calibrated for individual ESMs given fundamental inter-model differences. We then emulate 40 climate models of the Coupled Model Intercomparison Project Phase 5 (CMIP5) to create a ”super-ensemble”, i.e., a large ensemble ^[.7]which closely resembles a multi-model initial-condition ensemble. ^[.8]The thereby emerging ESM-specific emulator ^[.9]parameters provide essential insights on inter-model ^[.10]differences across a broad range of scales ^[.11]and characterize core properties of each ESM. Our results highlight that, for temperature at the spatio-temporal scales considered here, it is likely more advantageous to invest computational resources into generating multi-model ensembles rather than large single-model initial-condition en-

¹removed: spatially and temporally correlated

²removed: which closely resemble ESM runs spanning from 1870 to 2099.

³removed: framework that

⁴removed: emulator

⁵removed: mean temperature emulator

⁶removed: emulator

⁷removed: that

⁸removed: Furthermore, the

⁹removed: calibration

¹⁰removed: divergences

¹¹removed: which can be viewed as a ”model ID” of

sembles. Such multi-model ensembles can [..¹²] be extended to super-ensembles with [..¹³] emulators like the one presented here.

20 *Copyright statement.* TEXT

1 Introduction

The range of simulated climate responses to external radiative forcing is affected by both internal variability and inter-model differences (Hawkins and Sutton, 2009; Deser et al., 2012; Taylor et al., 2012). While inter-model uncertainty is typically accounted for by considering simulations from several climate models (Meehl et al., 2007; Taylor et al., 2012; Eyring et al., 2016), uncertainty due to internal climate variability is [..¹⁴] often quantified through running the same climate model a number of times with slightly different initial conditions [..¹⁵] (Deser et al., 2012; Fischer et al., 2013; Kay et al., 2015; Leduc et al., 2019).

As climate model ensembles are inherently expensive to run, there is an interest in approximating Earth System Model (ESM) output by computationally cheap emulators. In the field of climate science, the term emulator is used for a variety of statistical models which learn from existing runs of complex climate models to infer properties of runs which have not been generated yet. This [..¹⁶] makes it possible to explore the phase space at a lower computational cost. ESM emulators target different aspects of the climate system. For example, some emulators focus on the impacts of sub-grid scale parameterizations (Rougier et al., 2009; Williamson et al., 2013). Others target the effect of greenhouse gas emission scenarios on global mean temperature (Meinshausen et al., 2011; Goodwin, 2016) or on regional mean climate [..¹⁷] fields (Santer et al., 1990; Tebaldi and Arblaster, 2014; Tebaldi and Knutti, 2018). [..¹⁸] There are also emulators for regional-scale internal climate variability (Castruccio and Genton, 2016; Alexeeff et al., 2018; Link et al., 2019). Recently, first attempts have [..¹⁹] been made to emulate the full dynamics of simple general circulation models [..²⁰] (Scher, 2018; Scher and Messori, 2019).

In this study, the term emulator is used to refer to [..²¹] computationally cheap statistical [..²²] tools which generate additional realizations of [..²³] land temperature field time series for a specific greenhouse gas emission pathway at a yearly resolution. The presented emulator thus produces realizations which closely resemble initial-condition ensemble members of

¹²removed: then

¹³removed: geographically-explicit temperature

¹⁴removed: usually

¹⁵removed: (Deser et al., 2012; Fischer et al., 2013; Kay et al., 2015; Leduc et al., 2019; Maher et al., 2019)

¹⁶removed: allows for an exploration of the emulated

¹⁷removed: variable

¹⁸removed: Furthermore, there are emulators focusing on

¹⁹removed: additionally

²⁰removed: (Scher, 2018; ?).

²¹removed: a

²²removed: tool which generates

²³removed: geographically-explicit and spatio-temporally correlated

the considered ESMs. In the context of large multi-model ensembles, [..²⁴]our computationally cheap emulator can be used to produce look-alikes of large initial-condition ensembles for every model within the multi-model ensemble resulting in a ”super-ensemble”, i.e., a large ensemble which closely resembles a multi-model initial-condition ensemble.

To build this statistical temperature emulator, an overarching modular [..²⁵]framework is proposed and put into context of
45 previous work in Sect. 2. The employed data and terminology is described in Sect. 3, and the specific implementation of the [..²⁶]framework is introduced in Sect. 4. To visualize the characteristics and capabilities of the emulator, [..²⁷]detailed results are shown for four [..²⁸]example ESMs in Sect[..²⁹]. 5, before applying the emulator to the large CMIP5 (Coupled Model Intercomparison Project Phase 5, Taylor et al., 2012) multi-model ensemble containing 40 climate models in Sect. 6. [..³⁰]In Sect. 7, the results are discussed and finally, in Sect. 8, the conclusions and an outlook are provided.

50 2 A framework for end-to-end climate model emulation

We propose an additive framework for temperature emulation at the yearly scale for a specific greenhouse gas emission pathway which can be summarized as

$$T_{s,t} = f(T_t^{glob}) + \eta_{s,t}, \quad (1)$$

where the [..³¹]local temperature $T_{s,t}$ at grid point s and time t is defined as a [..³²]response to the [..³³]global mean temperature [..³⁴] T_t^{glob} , indicated by the function $f()$, and a stochastic local residual temperature variability term $\eta_{s,t}$. [..³⁵]Contributions from physical feedbacks other than the ones captured within the global mean temperature signal are thus neglected. The assumption of an underlying additivity is in line with frequently employed approaches in [..³⁶]uncertainty analysis in climate science (Hawkins and Sutton, 2009) and in climate change detection and attribution studies (Allen and Stott, 2003).

60 Our framework requires three [..³⁷]modules: a global mean temperature [..³⁸]module, a module for the grid-point level temperature response to the global mean temperature, and a local residual temperature variability [..³⁹]module. In the fol-

²⁴removed: the

²⁵removed: emulation

²⁶removed: proposed

²⁷removed: calibration results and example realizations

²⁸removed: selected

²⁹removed: 5.In Sect. 6, the emulator is applied to a large

³⁰removed: In addition to calibration results and example realizations, quantitative verification of in-sample and out-of-sample performance are presented.

³¹removed: temperature T

³²removed: deterministic

³³removed: current

³⁴removed: T_t^{glob}

³⁵removed: Deterministic contributions

³⁶removed: climate science uncertainty analysis (Hawkins and Sutton, 2009) and

³⁷removed: sub-modules. An emulator for

³⁸removed: , an emulator for the deterministic mean

³⁹removed: emulator

lowing, we place existing literature within these [..⁴⁰] modules before discussing the connections to our emulator. As this study is primarily concerned with temperature, we focus solely on this variable in our literature review. However, several of the referred studies treat also additional variables such as precipitation (e.g., Tebaldi and Arblaster, 2014; Seneviratne et al., 2016; Wartenburger et al., 2017) or cloud cover (e.g., Osborn et al., 2016).

[..⁴¹][..⁴²]

2.1 Global mean temperature module

Global mean temperature is often an output of computationally efficient simple energy-balance climate models (Meinshausen et al., 2011; Goodwin, 2016). While such models provide an estimate of the [..⁴³] global mean temperature trend, they do not produce interannual global mean temperature variability. To obtain an ensemble of global mean temperature variability, statistical models which account for temporal autocorrelation [..⁴⁴] can be used [..⁴⁵] (Brown et al., 2015).

[..⁴⁶]

2.2 Local temperature response module

Pattern scaling is a frequently employed approach to relate [..⁴⁷] local temperature to global mean temperature and is also used to emulate warming patterns across emission scenarios (Santer et al., 1990; Mitchell, 2003; Tebaldi and Arblaster, 2014). It was originally introduced by Santer et al. (1990) and different implementations exist (Mitchell, 2003). Most often, [..⁴⁸] temperature fields are averaged over a late 21st century multi-decadal time period and the associated average global mean temperature is obtained (Tebaldi and Arblaster, 2014). This pattern is then linearly interpolated to a desired global mean temperature [..⁴⁹]. An alternative is to extract the pattern [..⁵⁰] from a transient simulation at the time when [..⁵¹] the simulation reaches the desired global mean temperature [..⁵²] (Herger et al., 2015; Seneviratne et al., 2016; King et al., 2017). Other approaches include carrying out a linear regression (Lynch et al., 2017) or fitting a linear mixed-effect model (Alexeeff et al., 2018) to global mean temperature at each grid point individually.

The most important assumption underlying pattern scaling is that local mean temperatures are linearly related to global mean temperature and that this relationship is consistent across forcing scenarios. For surface temperature on land this assumption is satisfactorily met (Mitchell, 2003; Tebaldi and Arblaster, 2014; Seneviratne et al., 2016; Wartenburger et al., 2017; Osborn

⁴⁰removed: sub-modules

⁴¹removed: *Global mean temperature emulators*

⁴²removed: Emulated global

⁴³removed: forced global temperature response, they cannot quantify the effects of interannual climate variability. Statistical models accounting

⁴⁴removed: , however,

⁴⁵removed: to generate stochastic ensembles of unforced interannual variability of global mean temperature (Brown et al., 2015).

⁴⁶removed: *Local mean temperature emulators*

⁴⁷removed: grid-point level

⁴⁸removed: the temperature pattern provided by a climate model is

⁴⁹removed: level

⁵⁰removed: of change

⁵¹removed: this

⁵²removed: on its way to higher warming

et al., 2018). However, for strong mitigation scenarios and under strong aerosol forcing, pattern scaling is less accurate (May, 2012; Levy et al., 2013). [..⁵³]Additionally, it is assumed that external forcing and internal variability are independent which may not always be true [..⁵⁴](Lopez et al., 2014).

More complex local [..⁵⁵]response emulation methods are rare and often directly conditioned on CO₂ concentration profiles instead of global mean temperature (Castruccio et al., 2014; Holden and Edwards, 2010). For instance, it has been proposed to employ past trajectories of atmospheric CO₂ to model regional temperatures with an infinite distributed lag model to capture non-linear behaviour in spatial patterns for [..⁵⁶]regional-scale emulation (Castruccio et al., 2014) and within global space-time [..⁵⁷]models (Castruccio and Stein, 2013). Other authors use singular value decomposition to emulate decadal temperature fields across scenarios while accounting for complex spatio-temporal feedbacks (Holden and Edwards, 2010; Holden et al., 2014).

[..⁵⁸]While the focus is usually set on emulating the pattern associated with the global mean temperature trend, patterns associated with physical modes of variability such as the El Niño Southern Oscillation and the Pacific Decadal Oscillation can additionally be derived (McKinnon and Deser, 2018).

2.3 Local residual temperature variability module

Several approaches exist to emulate local residual temperature variability based on observations and climate model simulations [..⁵⁹](Castruccio and Stein, 2013; Osborn et al., 2016; McKinnon et al., 2017; Alexeeff et al., 2018; Link et al., 2019). Observations can be employed to avoid climate model biases but are limited [..⁶⁰]by rather short observational records when deriving the local temperature variability properties (Osborn et al., 2016; McKinnon et al., 2017; McKinnon and Deser, 2018). The simplest approach is to detrend observed [..⁶¹]temperature time series and obtain additional [..⁶²]realizations by shifting the starting date of the time series (Osborn et al., 2016). More realizations have been generated [..⁶³]by resampling spatial fields of detrended observed local temperature variability in blocks of two years [..⁶⁴](McKinnon et al., 2017). The approach was later refined to explicitly account for physical modes of variability to further reduce temporal autocorrelation in the resampled fields (McKinnon and Deser, 2018).

[..⁶⁵]

⁵³removed: An additional assumption of pattern scaling is

⁵⁴removed: and could lead to pattern estimation errors (Lopez et al., 2014).

⁵⁵removed: mean temperature

⁵⁶removed: regional scale emulation (Castruccio et al., 2014) as well as within a

⁵⁷removed: model

⁵⁸removed: *Local residual temperature variability emulators*

⁵⁹removed: (Osborn et al., 2016; McKinnon et al., 2017; Castruccio and Stein, 2013; Link et al., 2019; Alexeeff et al., 2018)

⁶⁰removed: to

⁶¹removed: time series of temperature

⁶²removed: realization

⁶³removed: from an observational record

⁶⁴removed: at the cost of losing longer-term memory (McKinnon et al., 2017; McKinnon and Deser, 2018)

⁶⁵removed: When employing ESM temperature time series

110 When employing ESMs instead, longer time series ^[..66] and multiple realizations are available to derive the statistical properties of the local residual temperature variability. Several authors fit autoregressive (AR) models to a set of climate model runs to account for temporal autocorrelation when ^[..67] emulating local residual temperature variability (Castruccio and Stein, 2013; Castruccio et al., 2014; Castruccio and Genton, 2016; Bao et al., 2016). ^[..68] Thereby, the spatial dependence in the innovation terms of the AR models can be considered by parameterizing their covariance structure with a Matérn covariance function (Castruccio and Stein, 2013; Castruccio and Genton, 2016; Bao et al., 2016). ^[..69] Alternatively, detrended ESM runs can be decomposed into their principal components ^[..70] and their phases can be randomly perturbed to generate additional realizations of local residual temperature variability (Link et al., 2019).

All approaches listed so far ^[..71] ^[..72] rely on the assumption that local residual temperature variability is stationary in time which is known not to be fulfilled everywhere. Olonscheck and Notz (2017) and references therein provide a comprehensive overview on possible changes in temperature variability in the historical time period and ^[..73] the business-as-usual greenhouse gas emission scenario for ^[..74] the large CMIP5 ^[..75] multi-model ensemble. They find that the strongest and most likely changes will occur over oceans but also point out land regions where variability is projected to change in the future. During the historical time period, they identify only weak changes in the variability. To account for such temporal non-stationarities, it has been proposed to resample detrended temperature fields of large single-model initial-condition ensembles within a certain window size around a global mean temperature level (Alexeeff et al., 2018). To enlarge the number of fields to sample from, a method has additionally been developed to stochastically emulate spatially non-stationary Gaussian fields with a LatticeKrig model (Nychka et al., 2018).

2.4 This study

130 While most studies focus on one or two of the modules required to mimic an initial-condition ensemble, this study proposes a framework which incorporates all three components. Since only 12 out of 40 CMIP5 models provide several initial-condition members, it is essential to test to what extent an emulator trained on a single run is able to approximate both its training run and additional independent initial-condition members. We thus emulate the full CMIP5 multi-model

⁶⁶removed: as well as multiple climate run

⁶⁷removed: simulating additional realizations of

⁶⁸removed: Most of them furthermore consider

⁶⁹removed: An alternative approach is based on decomposing

⁷⁰removed: , whose phases are

⁷¹removed: assume temporal stationarity in the local residual climate variability. Large single-model initial-condition ensembles, however, can be used to re-sample detrended temperature fields within a certain window size around a global mean temperature level which makes it possible to account for temporal non-stationarities in local residual temperature variability in a warming world (Alexeeff et al., 2018). To enlarge the amount of fields to sample from, a method has additionally been proposed to stochastically simulate spatially non-stationary Gaussian fields with a LatticeKrig model (Nychka et al., 2018).

⁷²removed: Overall, most studies

⁷³removed: a

⁷⁴removed: models in the large Coupled Model Intercomparison Project Phase 5 (

⁷⁵removed:)

Table 1. Terms used to refer to different climate model runs throughout this study.

Name	Description	Application	Figures
Training run	Climate model run (1870–2100) used to calibrate the emulator parameters	- emulator calibration - emulator evaluation in terms of fitting the training run	3, 7, 8, 9, 10, 11, 12
Test run	Independent initial-condition ensemble [⁸²]member (1870–2100) not used to calibrate the emulator parameters	- emulator evaluation in terms of mimicking a climate model initial-condition ensemble	4, 5, 6, 8, 9, 10, 11, 12

ensemble based on single training runs and create a super-ensemble which accounts for inter-model uncertainty across
 135 all 40 climate models. To the best of our knowledge, this study is the first to implement an emulator which mimics an
 initial-condition ensemble based on a single training run and applies it to such a large multi-model ensemble.

3 Data and terminology

3.1 Data sources and terminology

Runs from 40 CMIP5 climate models (Taylor et al., 2012) covering the historical time period (1870–2005) and the business-
 140 as-usual greenhouse gas emission scenario RCP8.5 [⁷⁶](2005–2100, Riahi et al., 2011) are employed. [⁷⁷]To calibrate
 the emulator, a single run per climate model is used. This run is referred to as the training run (Table 1). For 12 out of
 40 CMIP5 climate models more than one initial-condition member [⁷⁸]is available. These additional independent initial-
 condition ensemble members are referred to as test runs (Table 1). A special focus is set on four ESMs [⁷⁹]with differing
 model genealogies (Knutti et al., 2013), namely CanESM2, CESM1(CAM5), [⁸⁰]HadGEM2-ES, and MPI-ESM-LR[⁸¹
 145]. All climate models, the associated modeling groups, and the number of initial-condition members employed here are
 listed in Table A1.

Additionally, stratospheric aerosol optical depth is used as a proxy for volcanic activity during the historical time period.
 This aerosol dataset was originally described by Sato et al. (1993) and later updated to cover the considered time period.

3.2 Data processing

150 Here, we focus on surface temperature anomaly at a yearly resolution. Temperature fields were bilinearly interpolated onto a
 2.5° x 2.5° grid resulting in 3043 land grid points for each climate model. Yearly mean temperatures were computed at each grid
 point and the average over the reference period of 1870–1899 in the training run at the respective grid points was subtracted.

⁷⁶removed: (2005–2099, Riahi et al., 2011)

⁷⁷removed: The models, the associated modeling groups, and the number of

⁷⁸removed: used here are listed in Table A1. Out of these 40 models, a

⁷⁹removed: (

⁸⁰removed: HadGEM2-ESM,

⁸¹removed:). These are chosen because they represent different model genealogies according to Knutti et al. (2013) and because they provide several

In the text, for simplicity reasons, we use the term "temperature" when referring to "yearly surface temperature anomaly". For the stratospheric aerosol optical depth, the globally averaged yearly time series is employed.

155 Whenever regional averages are shown, area-weighted means are referred to. The regions employed in this study are 26 SREX land regions (Seneviratne et al., 2012) as well as global mean and global land mean (Fig. 1). While global mean refers to the average across all grid points, global land mean refers to the average across all land grid points excluding Antarctica.

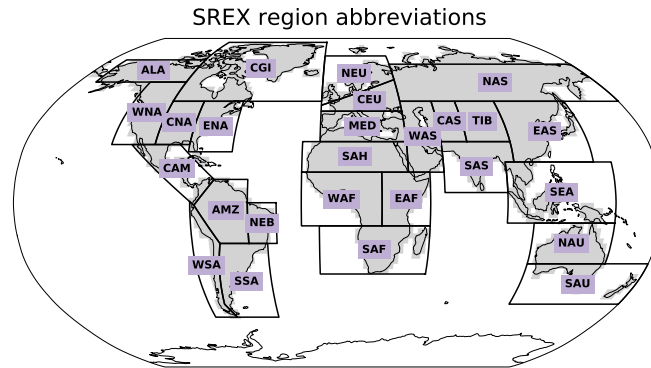


Figure 1. Map of the SREX regions and their abbreviations. The considered land grid points are shown in grey.

4 Methods

160 4.1 [\[..⁸³\]](#) Framework implementation

4.1.1 General approach

[\[..⁸⁶\]](#) We follow the framework introduced in Sect. 2 to emulate temperature fields at the yearly scale for a specific greenhouse gas emission pathway [\[..⁸⁷\]](#)

[\[..⁸⁸\]](#)

165 [\[..⁸⁹\]](#). The chosen implementation is shown in Fig. 2 and detailed information for each individual module is provided in the following sections. In short, the global [\[..⁹⁰\]](#) mean temperature T_t^{glob} is split into a [\[..⁹¹\]](#) trend and a variability term,

⁸³removed: General approach and terminology

⁸⁶removed: To

⁸⁷removed: , we follow the framework introduced in Sect. 2. Figure 2 shows the chosen implementation graphically and the associated governing equation reads

⁸⁹removed: Hence

⁹⁰removed: forcing

⁹¹removed: deterministic trend $T_t^{glob,det}$ and a stochastic variability $T_t^{glob,var}$

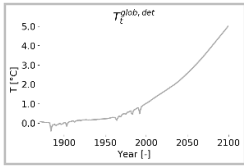
$$\text{Emulator: } T_{s,t} = f(T_t^{glob}) + \eta_{s,t}$$

Global mean temperature T_t^{glob} module

$$T_t^{glob, sm} + T_t^{glob, volc} + \alpha_0 + \sum_{k=1}^{k=p} \alpha_k \cdot T_{t-k}^{glob, var} + \varepsilon_t, \quad \varepsilon_t \sim \mathcal{N}(0, \sigma)$$

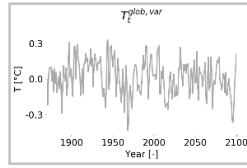
Global mean temperature trend

$T_t^{glob, trend}$



Global mean temperature variability

$T_t^{glob, var}$



Local response $T_{s,t}^{resp}$ module

$$\beta_S^{trend} \cdot T_t^{glob, trend} + \beta_S^{int} + \beta_S^{var} \cdot T_t^{glob, var}$$

Local residual variability $\eta_{s,t}$ module

$$\gamma_{0,s} + \gamma_{1,s} \cdot \eta_{s,t-1} + \nu_{s,t} =$$

$$\nu_{s,t} \sim \mathcal{N}(0, \Sigma_\nu(r))$$

$T_{s,t}$

Emulation

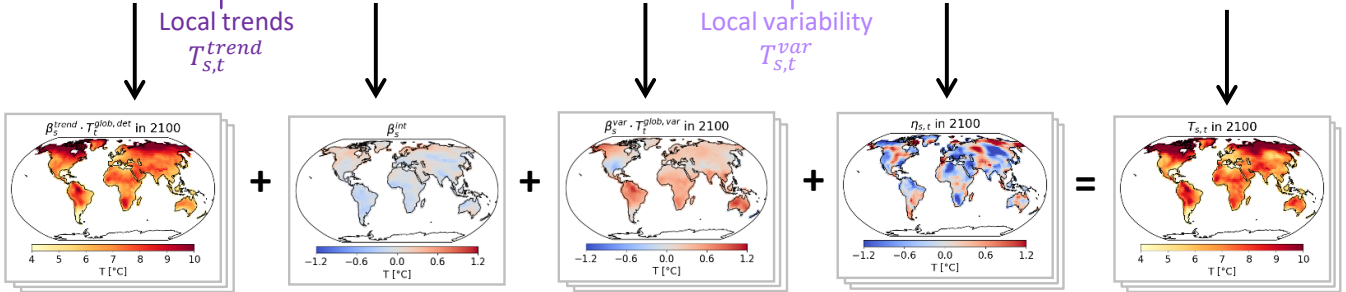


Figure 2. Illustration of the [..⁸⁴] emulation framework.[..⁸⁵]

both of which contribute linearly to the local [..⁹²] temperature $T_{s,t}$ [..⁹³]. The residual [..⁹⁴] local temperature variability $\eta_{s,t}$ is modeled [..⁹⁵] as an AR(1) process with spatially correlated innovations.

⁹²removed: grid-point

⁹³removed: in a deterministic way via the function $f()$

⁹⁴removed: spatio-temporally correlated

⁹⁵removed: stochastically

[..⁹⁶] To calibrate the emulator, a single run spanning 231 years ([..⁹⁷] 1870–2100) per model is used. For the calibration, 170 the global mean temperature trajectory and the associated land temperature fields are required.

[..⁹⁸]

[..⁹⁹]

4.2 [..¹⁰⁰]

4.1.1 [..¹⁰¹]

175 [..¹⁰²]

4.1.1 Global mean temperature module

In the global mean temperature [..¹⁰³] module, additional realizations of global mean temperature time series [..¹⁰⁴] T_t^{glob} are generated. For this purpose, T_t^{glob} is separated into a [..¹⁰⁵] trend $T_t^{glob,trend}$ shared by all emulations and a [..¹⁰⁶] variability term $T_t^{glob,var}$ [..¹⁰⁷] which varies between individual emulations:

$$180 \quad T_t^{glob} = T_{[..¹⁰⁸]glob,trend}_t + T_t^{glob,var}. \quad (2)$$

In [..¹⁰⁹] $T_t^{glob,trend}$, smooth forcing $T_t^{glob,sm}$ and abrupt changes induced by volcanic eruptions $T_t^{glob,volc}$ are accounted for in an additive way:

$$T_{[..¹¹⁰]glob,trend}_t = T_t^{glob,sm} + T_t^{glob,volc}. \quad (3)$$

First, $T_t^{glob,sm}$ is derived by locally weighted scatterplot smoothing (LOWESS) of T_t^{glob} .

⁹⁶removed: For each considered climate model, an emulator is trained on a single climate model run spanning 230

⁹⁷removed: 1870–2099). To calibrate the emulators, the land temperature field time series as well as the

⁹⁸removed: The emulators' performance is subsequently evaluated with in-sample and, where possible, additionally, out-of-sample verification. Thereby, the in-sample verification is carried out on the training run itself and indicates how successfully our framework implementation captures the training run. For climate models with several initial-condition ensemble members, out-of-sample verification is conducted on the runs not employed during the training of the emulators. Hence, the out-of-sample verification serves as a proxy for the emulators' capability to mimic true ESM initial-condition ensembles.

⁹⁹removed: In the following sections, the chosen equations for each emulator module are introduced and the respective calibration procedures are described.

Afterwards, the emulator performance evaluation approach is explained.

¹⁰⁰removed: Framework implementation

¹⁰¹removed: Global mean temperature emulator

¹⁰²removed: The

¹⁰³removed: emulator stochastically creates

¹⁰⁴removed: . Global mean temperature

¹⁰⁵removed: deterministic part $T_t^{glob,det}$

¹⁰⁶removed: stochastic global variability part

¹⁰⁷removed: varying between realizations:

¹⁰⁹removed: the deterministic trend, $T_t^{glob,det}$,

185 In a next step, $[..^{111}]T_t^{glob,volc}$ is approximated as the linear response of the residuals of the smooth trend, i.e., $[..^{112}]T_t^{glob} - T_t^{glob,sm}$, to stratospheric aerosol optical depth AOD_t with regression coefficients λ_0 and λ_1 :

$$T_t^{glob,volc} = \lambda_0 + \lambda_1 \cdot AOD_t \quad (4)$$

The time series of $[..^{113}]$ global mean temperature variability $T_t^{glob,var} = T_t^{glob} - T_t^{glob,trend}$ is modeled as an AR process of order p with coefficients $\alpha_0, \dots, \alpha_p$ such that

$$190 \quad T_t^{glob,var} = \alpha_0 + \sum_{k=1}^{k=p} \alpha_k \cdot T_{t-k}^{glob,var} + \epsilon_t \quad \text{with} \quad \epsilon_t \sim \mathcal{N}(0, \sigma), \quad (5)$$

whereby ϵ_t is a white noise innovation term drawn from a Gaussian distribution with mean zero and $[..^{114}]$ standard deviation $\sigma[..^{115}]$.

$[..^{116}]$ In this study, the LOWESS smoothing window length is 50 years with weights decaying with increasing distance according to a tricube weight function. The regression coefficients for the forced response to volcanic eruptions are obtained with the ordinary least squares (OLS) algorithm. The coefficients of the AR process are fit by means of maximum likelihood and the Bayesian Information Criterion (BIC) $[..^{117}]$ is employed to select $[..^{118}]$ its order p with the maximum $[..^{119}]$ considered order being eight.

4.1.2 Local $[..^{120}]$ temperature $[..^{121}]$ $[..^{122}]$ response module

The local temperature response module translates the global mean temperature signal into a grid-point level response $T_{s,t}^{resp}$. $[..^{123}]$ Motivated by the pronounced linear scaling of regional land temperatures with global mean temperature (Seneviratne et al., 2016; Wartenburger et al., 2017), the local $[..^{124}]$ response is expressed as

$$T_{s,t}^{resp} = f(T[..^{125}]_t^{glob}) = f(T_s^{glob,trend}, T_t^{glob,var}) = \beta[..^{126}]_s^{trend} \cdot T[..^{127}]_t^{glob,trend} + \beta_s^{int} + \beta_s^{var} \cdot T_t^{glob,var} [..^{128}], \quad (6)$$

with regression coefficients $[..^{129}]$ β_s^{trend} , β_s^{int} , and β_s^{var} whereby β_s^{int} represents the intercept term. Hence, the response of the local mean temperature to $[..^{130}]T_t^{glob,trend}$ and $T_t^{glob,var}$ are separately taken into account.

¹¹¹removed: $T_t^{glob,volc}$

¹¹²removed: $(T_t^{glob} - T_t^{glob,sm})$

¹¹³removed: unforced global temperature variability $T_t^{glob,var} = T_t^{glob} - T_t^{glob,det}$

¹¹⁴removed: its standard deviation set to the empirical standard deviation of the innovations of the training samples (

¹¹⁵removed:).

¹¹⁶removed: Here

¹¹⁷removed: , which punishes model complexity to avoid overfitting,

¹¹⁸removed: the

¹¹⁹removed: possible order set to 8.

¹²⁰removed: mean

¹²¹removed: emulator

¹²²removed: The local mean temperature emulator

¹²³removed: Due to

¹²⁴removed: mean temperature $T_{s,t}^{resp}$

¹²⁹removed: β_s^{det} , β_s^{var} , and

¹³⁰removed: the deterministic global trend $T_t^{glob,det}$ and the response to the global temperature variability

205 In this study, the linear regression coefficients are estimated with OLS at each grid point.

4.1.3 Local residual temperature variability [..¹³¹]module

The local residual temperature variability $\eta_{s,t}$ refers to the spatio-temporally correlated residual variability which cannot be accounted for through a response to [..¹³²] T_t^{glob} . This variability is assumed to be [..¹³³] Gaussian in nature (see S1 for the results of a Shapiro-Wilks test for normality) and stationary in time which makes it possible to model the time series as local
 210 AR(1) processes with spatially correlated innovations (Humphrey and Gudmundsson, 2019). Hence, additional realizations of [..¹³⁴] $\eta_{s,t}$ are generated stochastically according to

$$\eta_{s,t} = \gamma_{0,s} + [..¹³⁵]\gamma[..¹³⁶] $_{1,s} \cdot \eta[..¹³⁷] $_{s,t-1} + \nu_{s,t} \quad \text{with} \quad \nu_{s,t} \sim \mathcal{N}(0, \Sigma_\nu(r)), \quad (7)$$$$

whereby $\gamma_{0,s}$ [..¹³⁸] and $\gamma_{1,s}$ are the coefficients of the AR model and $\nu_{s,t}$ are [..¹³⁹] spatially correlated innovations drawn from a multivariate Gaussian with mean zero and [..¹⁴⁰] covariance matrix $\Sigma_\nu(r)$ (Cressie and Wikle, 2011).

215 For an AR(1) process, $\Sigma_\nu(r)$ can be analytically derived from the covariance matrix of the residual variability $\Sigma_\eta(r)$ with

$$\Sigma_\nu(r)_{i,j} = \sqrt{1 - \gamma_{1,i}} \cdot \sqrt{1 - \gamma_{1,j}} \cdot \Sigma_\eta(r)_{i,j}, \quad (8)$$

whereby the indices i and j refer to spatial locations s (Cressie and Wikle, 2011).

[..¹⁴¹] To estimate $\Sigma_\eta(r)$, the empirical covariance matrix [..¹⁴²] $\tilde{\Sigma}_\eta$ is computed. However, $\tilde{\Sigma}_\eta$ is rank deficient because
 220 substantially fewer temperature field samples are available than there are land grid points. Thus, $\tilde{\Sigma}_\eta$ needs to be regularized to obtain a robust estimate of the co-variations between the grid points. For this purpose, [..¹⁴³] we employ localization, an approach which is well established in the field of data assimilation (Carrassi et al., 2018). Localization retains anisotropy on regional scales which is an important asset when stochastically modeling temperature variability since anisotropy is a prevalent feature due to physical factors such as prevailing wind direction and geometry of mountainous terrain. To
 225 localize $\tilde{\Sigma}_\eta$, it is point-wise multiplied with a smooth correlation function $G(r)$ with exponentially vanishing correlations with distance:

$$\Sigma_\eta(r) = \tilde{\Sigma}_\eta \circ G(r), \quad (9)$$

¹³¹removed: emulator

¹³²removed: the global mean temperature signal

¹³³removed: stationary in time and

¹³⁴removed: the local residual temperature variability are generated according to

¹³⁸removed: ..., $\gamma_{p_s, s}$ refer to

¹³⁹removed: the white noise

¹⁴⁰removed: the regularized empirical spatial covariance matrix $\Sigma(r)$

¹⁴¹removed: The rank deficient

¹⁴²removed: $\tilde{\Sigma}$

¹⁴³removed: $\tilde{\Sigma}$ is localized by multiplying it

whereby \circ denotes the Hadamard product. Here, $G(r)$ is the numerically efficient Gaspari-Cohn function (Gaspari and Cohn, 1999) which vanishes beyond two times the localization radius L :

$$230 \quad G(r) = \begin{cases} 1 - \frac{5}{3} \cdot r^2 + \frac{5}{8} \cdot r^3 + \frac{1}{2} \cdot r^4 - \frac{1}{4} \cdot r^5, & \text{if } 0 \leq r < 1, \\ 4 - 5 \cdot r + \frac{5}{3} \cdot r^2 + \frac{5}{8} \cdot r^3 - \frac{1}{2} \cdot r^4 + \frac{1}{12} \cdot r^5 - \frac{2}{3} \cdot r^{-1}, & \text{if } 1 \leq r < 2, \\ 0, & \text{if } r \geq 2, \end{cases} \quad (10)$$

with $r = \frac{d}{L}$ and d the geographical distance between two grid points.

[..¹⁴⁴]

In this study, the AR(1) coefficients are fit at each grid point by means of maximum likelihood. [¹⁴⁵] In our framework implementation, the obtained intercept terms $\gamma_{0,s}$ are effectively zero, as the local response module already contains an intercept term (Eq.6). The localization radius to regularize [¹⁴⁶] $\tilde{\Sigma}_\eta$ is determined by cross-validation [¹⁴⁷] with a leave-one-out approach. Localization radii between 1000 and 4750 km every 250 km [¹⁴⁸] are tested. Thereby, the empirical covariance matrix is estimated based on 230 years and the likelihood to draw the [¹⁴⁹] field of the left-out year from the regularized matrix is computed. [¹⁵⁰] This process is repeated until every year has been left out once for every localization radius. The respective log-likelihood values for each localization radius are summed up [¹⁵¹] across the left-out years and
 240 the radius which is associated with the maximum likelihood is chosen.

4.2 Evaluating the emulator

[¹⁵²] The emulator's performance is evaluated on the training run and - where available - on test runs. While the evaluation on the training run indicates how successfully this framework implementation captures the training run, the evaluation on the test runs serves as a proxy for the emulator's capability in mimicking true ESM initial-condition ensembles. For the
 245 evaluation, 1000 emulations are generated for each climate model.

[..¹⁵³]

¹⁴⁴removed: An ARprocess is

¹⁴⁵removed: BIC is employed to select the order of the AR processes with the maximum order set to 4. In a subsequent step, the

¹⁴⁶removed: the empirical covariance matrix of the innovations with

¹⁴⁷removed: . For this purpose, the time series of 230 years is split into 5 folds. Each fold contains 1 block of 6 years and 8 blocks of 5 years spread out evenly across the ESM run with maximum spacing in time between individual blocks. The empirical covariance matrix of the innovations is estimated based on 4 folds and regularized with localization

¹⁴⁸removed: . The

¹⁴⁹removed: innovations of the 5th fold from each one of these regularized covariance matrices

¹⁵⁰removed: The likelihood values

¹⁵¹removed: over all folds and finally, the localization radius

¹⁵²removed: To evaluate the emulator, the focus is set on ensemble reliability, i. e., the ability to capture the distribution of ESM runs with an ensemble of emulations (Weigel, 2012). This approach is chosen because the generated emulations cannot be directly compared to ESM runs as they differ in each realization and thus the traditionally employed measure of "best-fit", i.e. , least deviation from the ESM run, is not a meaningful metric.

¹⁵³removed: Visual

4.2.1 Local trends verification

The local trends $T_{s,t}^{trend}$ are shared by all emulations and serve as an estimate of the externally forced response with

$$T_{s,t}^{trend} = \beta_s^{trend} \cdot T_t^{glob,trend} + \beta_s^{int}. \quad (11)$$

250 To evaluate how well the emulated local trends capture true climate model runs, the Pearson correlation of $T_{s,t}^{trend}$ with $T_{s,t}$ of the corresponding training run is computed. For climate models with test runs, the correlation coefficient is additionally computed between $T_{s,t}^{trend}$ and each test run.

4.2.2 Local variability verification

The local variability $T_{s,t}^{var}$ is different in each emulation and corresponds to the internally generated natural variability:

$$255 T_{s,t}^{var} = \beta_s^{var} \cdot T_t^{glob,var} + \eta_{s,t}. \quad (12)$$

To compare the emulated $T_{s,t}^{var}$ to true climate model runs, an estimate for the local variability within the climate models needs to be obtained. For this purpose, the emulated local trends $T_{s,t}^{trend}$ (Eq. 11) are subtracted from the climate model $T_{s,t}$.

To evaluate $T_{s,t}^{var}$, on the grid-point level, lag-1 temporal autocorrelations and standard deviations are considered. 260 Additionally, spatial cross-correlations between grid points are verified. These quantities are computed for each individual emulation as well as [..¹⁵⁴] for all climate model runs. For each quantity, the Pearson correlation coefficient between each individual emulation and the training run is calculated. Additionally, the correlation between each individual test run and the respective training run is computed where test runs are available. These correlations between the climate model runs serve as benchmark values for the correlations between the emulations and the training run.

265 4.2.3 Regional-scale ensemble reliability verification

On regional scales, the emulated temperatures $T_{s,t}$ (Eq. 1) are evaluated visually and quantitatively in terms of [..¹⁵⁵] ensemble reliability, i.e., the ability to capture the distribution of ESM runs with an ensemble of emulations (Weigel, 2012). For the visual verification, regionally averaged emulated time series are compared to climate model runs for global land, Central Europe (CEU), and Southern South America (SSA)[..¹⁵⁶]. In the quantitative verification, the [..¹⁵⁷] emulator's 270 ability to reliably reproduce a set of ESM quantiles (5 %, 50 %, 95 %) [..¹⁵⁸] is evaluated in all 27 land regions. Smooth time

¹⁵⁴removed: quantitative reliability verification are carried out. The visual verification consists of a comparison between ESM runs and emulations

¹⁵⁵removed: temperature field snapshots and regionally averaged temperature time series . The time series are always shown for three specific regions (

¹⁵⁶removed:) which differ substantially in their temperature properties and thus highlight the emulator's flexibility in adapting to regional climate characteristics

¹⁵⁷removed: emulations'

¹⁵⁸removed: in the 26 considered SREX regions and global land is evaluated . For this purpose, an ensemble of 200 emulations is generated for each considered climate model

series of the emulated quantiles are obtained based on the 1000 emulations and the percentage of time slots during which ^[..159] a climate model run is below these emulated quantiles is counted. ^[..160]

^[..161] This is done for the training run and – where available – also for the test runs. Additionally, ^[..162] the counting is carried out for ^[..163] each individual emulation. The resulting deviations of the individual emulations from the emulated quantiles can be compared to the deviation the climate model runs exhibit from the emulated quantiles. If the climate model run deviation lies within the 95% interval spanned by the individual emulation deviations, the climate model run is considered indistinguishable from individual emulations at this quantile.

5 Exploring emulator properties for ^[..164] four example ESMs

5.1 Calibration results

280 The ^[..176] parameters obtained from training ^[..177] the emulator on four example ESMs reveal distinct inter-ESM differences in every emulator module (Fig. 3). The ^[..178] global mean temperature trends ^[..179] diverge by 0.9 °C by the end of the 21st century. For each ESM, ^[..180] $T_t^{glob,var}$ is described by oscillating AR coefficients with the first lag being positive, but the AR process order ^[..181] and the standard deviations of the innovations vary ^[..182].

In the local ^[..183] response module (Eq. 6), the strongest warming rates, i.e., the largest ^[..184] β_s^{trend} terms, are found in 285 the northern high latitudes, but there are substantial differences in the ^[..185] β_s^{trend} patterns between emulators trained on different ESMs (Fig. 3). For example, the CESM1(CAM5) emulator ^[..186] exhibits less warming in the tropics than the others ^[..187] do. In all emulators, the intercept term β_s^{int} is generally small in magnitude and smooth in space. The β_s^{var} fields

¹⁵⁹ removed: the regional temperature in an ESM run is higher than the respective quantile in the ensemble of emulations

¹⁶⁰ removed: The difference between the emulated and the counted quantile then gives an indication of the reliability of the emulated ensemble with respect to the ESM run at hand.

¹⁶¹ removed: In-sample verification on the training run is conducted for each one of the 40 CMIP5 models

¹⁶² removed: out-of-sample verification on ESM runs not seen during training

¹⁶³ removed: the 12 CMIP5 models with several initial-condition ensemble members. In the latter case, the differences between the emulated and counted quantiles are computed individually for each ESM initial-condition ensemble member not seen during training and are subsequently averaged.

¹⁶⁴ removed: selected

¹⁷⁶ removed: calibrated

¹⁷⁷ removed: an emulator on each one of the four selected

¹⁷⁸ removed: deterministic

¹⁷⁹ removed: $T_t^{glob,det}$ diverge by about 1

¹⁸⁰ removed: the global temperature variability

¹⁸¹ removed: p

¹⁸² removed: between them

¹⁸³ removed: mean

¹⁸⁴ removed: β_s^{det}

¹⁸⁵ removed: spatial patterns of β_s^{det}

¹⁸⁶ removed: warms less

¹⁸⁷ removed: .

Emulator calibration parameters

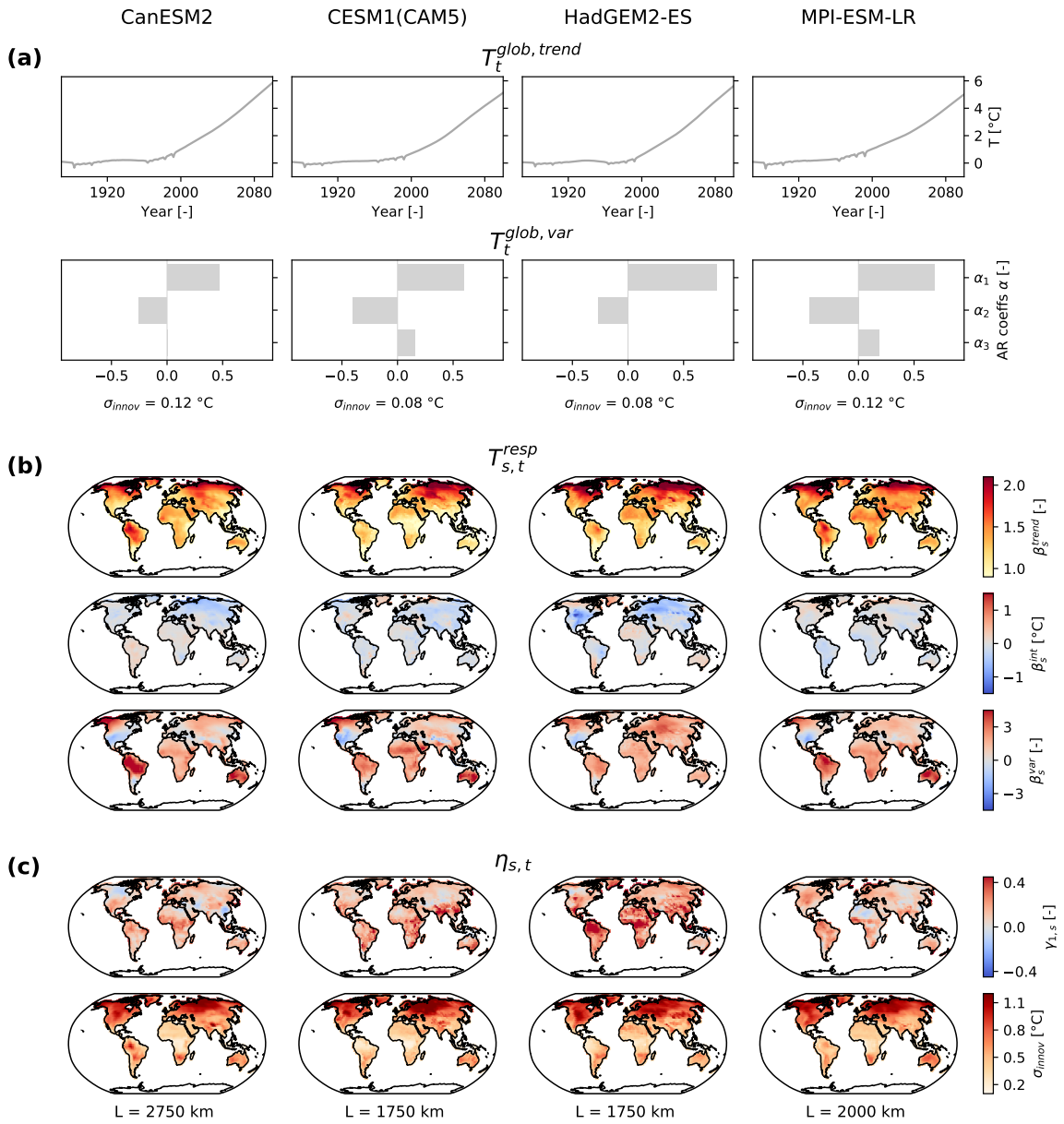


Figure 3. Emulator calibration parameters (rows) for [..¹⁶⁵] four [..¹⁶⁶] example ESMS (columns). (a) For the global mean temperature [..¹⁶⁷] module [..¹⁶⁸] $T_t^{glob, trend}$ and the AR coefficients plus the standard deviation of the innovations of [..¹⁶⁹] $T_t^{glob, var}$ are depicted. (b) For the local [..¹⁷⁰] temperature response module, the regression coefficients are shown. (c) For the local residual temperature variability [..¹⁷¹] module, the lag-1 AR [..¹⁷²] coefficients, the standard [..¹⁷³] deviations of the innovations [..¹⁷⁴], and the localization [..¹⁷⁵] radii are displayed.

indicate that Alaska, Amazon, ^[..188] and Australia frequently co-vary with global ^[..189] mean temperature variability. Only for HadGEM2-ES ^[..190] Central Asia emerges as a region of large β_s^{var} values. ^[..191]

290 ^[..192] The local residual variability (Eq. 7) ^[..193] exhibits generally less memory in the northern high latitudes than in the tropics ^[..194] as indicated by the lag-1 autocorrelation coefficients (Fig. 3). The innovations are largest in magnitude in high latitude continental climates such as North Asia and smallest in the tropics. However, also for these quantities the ^[..195] patterns differ between emulators calibrated on different ESMs. The localization radii chosen to regularize the empirical spatial covariance matrix ^[..196] $\tilde{\Sigma}_\eta$ range from 1750 to 2750 km.

295 5.2 Example realizations

Emulated temperature fields are visually indistinguishable from ESM test runs that were not used during training (Fig. 4). All fields exhibit the strongest warming ^[..197] and variability in the northern high latitudes. In terms of variability, CESM1(CAM5), HadGEM2-ES, and their emulations, show more patchy behaviour, i.e., locally more confined variability, than CanESM2 and MPI-ESM-LR.

300 Time series of emulations and ESM ^[..202] test runs averaged over global land, CEU, and SSA highlight the emulators capability to reproduce regionally characteristic behaviour of the climate system (Fig. 5). These regions differ in terms of ^[..203] warming trend and variability around this trend. The variability ^[..204] is smallest on the global scale since local anomalies tend to average out globally. In CEU, the warming rate as well as the variability ^[..205] are larger than in SSA.

5.3 Emulator transferability between ESMs

305 Figure 6 shows explicitly what the results of Sects. 5.1 and 5.2 have already hinted at, namely that an the ensemble of emulations generated by an emulator calibrated on a specific ESM is capturing unique properties of that ESM, which in turn are not transferable to other ESMs. For example, the warming rate of the ensemble generated by the CESM1(CAM5) emulator is inconsistent with ^[..211] all three other ESMs on the global land scale. As expected, differences are also found in the variability

¹⁸⁸ removed: Africa,

¹⁸⁹ removed: variability $T_t^{glob,var}$. Only in the

¹⁹⁰ removed: emulator

¹⁹¹ removed: In each one of the calibrated emulators, the intercept term β_s^{int} is generally small in magnitude and smooth in space hinting at overall successful fits.

¹⁹² removed: To model local

¹⁹³ removed: , AR processes of order 1 or even 0 suffice to capture the characteristic behaviour of each ESM in large parts of the globe (Fig. 3). Generally, there is

¹⁹⁴ removed: . The innovations of the local AR processes

¹⁹⁵ removed: spatial

¹⁹⁶ removed: of the innovations vary between 1500 and 2250

¹⁹⁷ removed: as well as

²⁰² removed: runs not employed during training

²⁰³ removed: underlying

²⁰⁴ removed: term

²⁰⁵ removed: around the trend

²¹¹ removed: ESM runs from

Temperature fields of ESM test runs and emulations in 2100

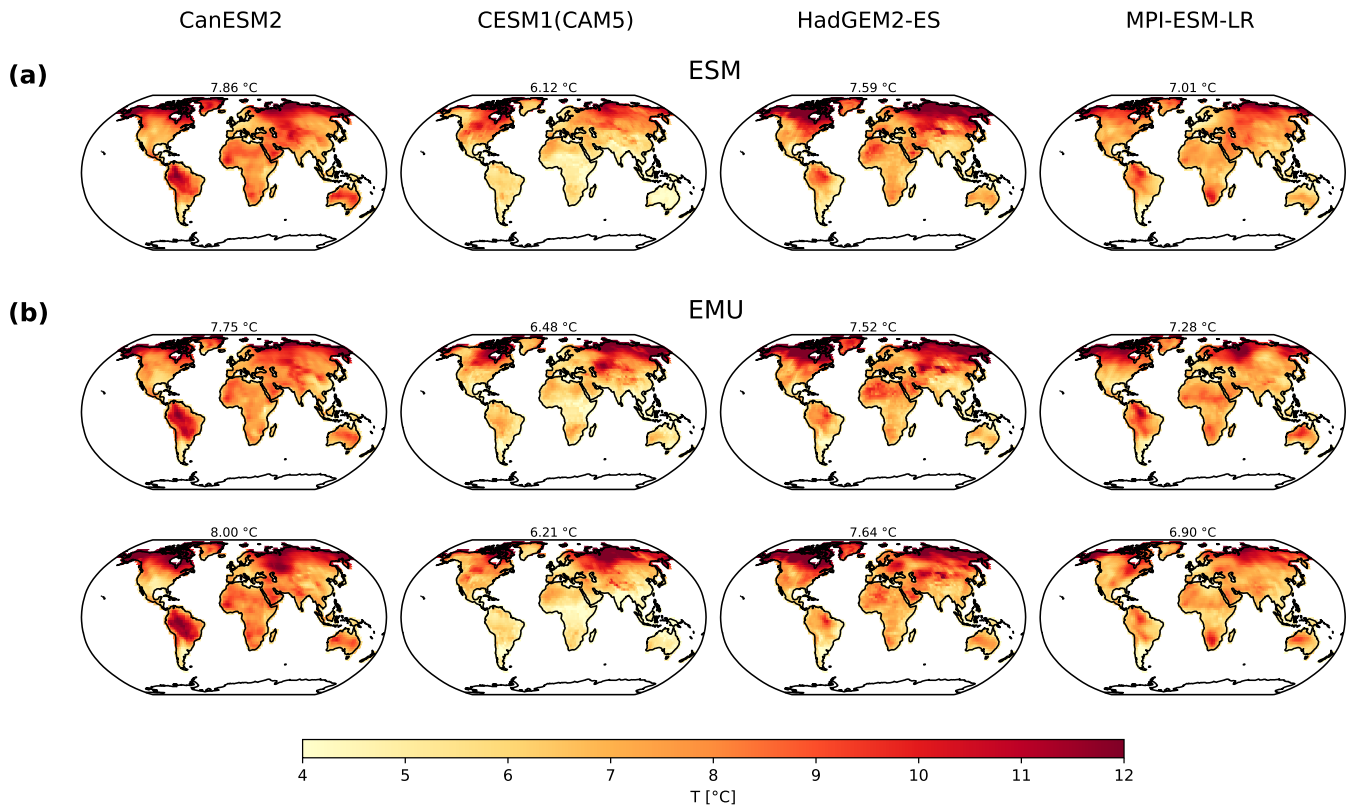


Figure 4. Temporal snapshots depicting temperature field realizations in [..¹⁹⁸]2100 (rows) for [..¹⁹⁹]four example ESMs (columns). (a) One ESM field [..²⁰⁰]from a test run and (b) two [..²⁰¹]emulations (EMUs) are shown. The temperature on top of each map refers to the global land mean.

around the trend which is, e.g., visibly smaller in SSA in the CESM1(CAM5) emulations than in the runs of the other ESMs.
 310 The implications of these results are further discussed in Sect. 7.3.

6 Creating a CMIP5 super-ensemble

[..²¹⁴]

²¹⁴removed: Based on the insights gained in Sect. 5, the 40-model CMIP5 ensemble is emulated by training an emulator for each one of the climate models.

Regionally averaged temperature time series of emulations and ESM test runs

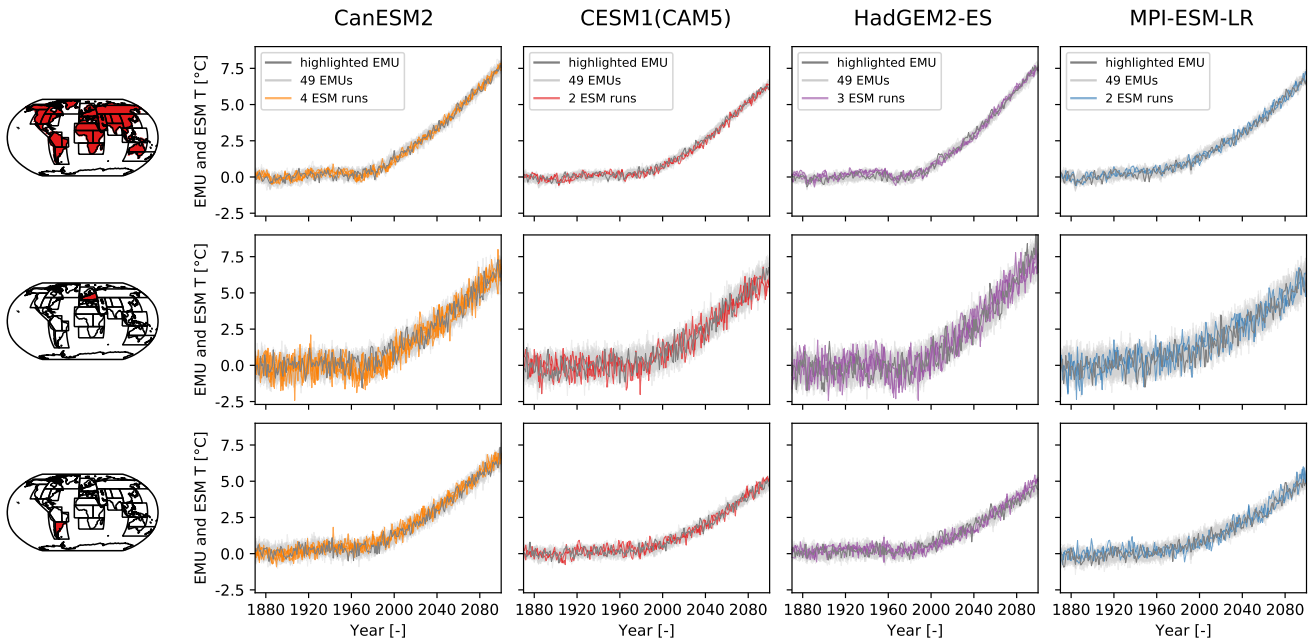


Figure 5. Regionally averaged temperature time series (rows) for [..²⁰⁶]four [..²⁰⁷]example ESMs (columns). The regions are from top to bottom: global land, Central Europe (CEU), and Southern South America (SSA). In each panel, 1 emulation (EMU) is highlighted in dark grey and [..²⁰⁸]49 other emulations are shown in light grey. Additionally, all available ESM [..²⁰⁹]test runs are plotted in color[..²¹⁰].

Inter-ESM transferability of the CESM1(CAM5) emulator in different regions

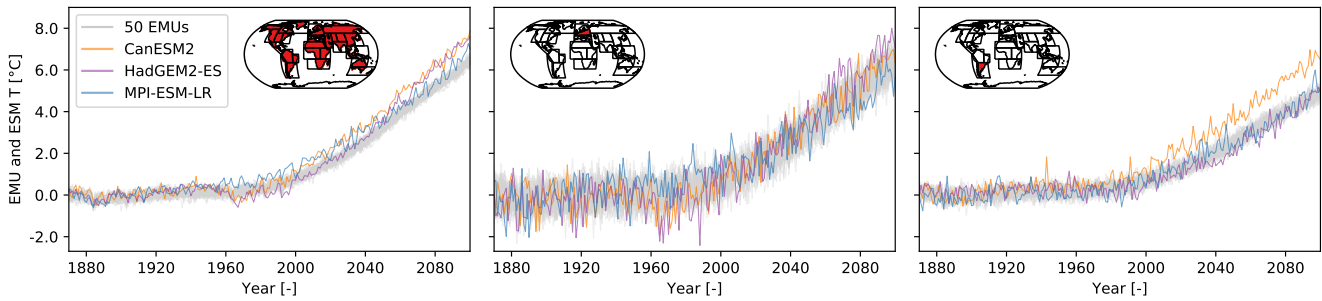


Figure 6. Time series of 50 emulations (EMUs) from the CESM1(CAM5) emulator (light grey) overlaid with runs from the three other [..²¹²]example ESMs [..²¹³]for three regions from left to right: global land, CEU, and SSA.

6.1 Calibration results

Figure 7 shows summary statistics of the calibrated parameters for each CMIP5 climate model highlighting inter-model differences in each emulator module. ^[..²¹⁵] In the supplementary information, plots analogous to Fig. 3 are **additionally** provided for each climate model for readers interested in the geographical patterns of the ^[..²¹⁶] emulator parameters (Figs. ^[..²¹⁷] S2–S10). ^[..²³³] In the global mean temperature ^[..²³⁴] module (Eq. 2), $T_t^{glob,trend}$ ranges between 3.4 and 6.3 °C at the end of the 21st century (Fig. 7). For 45 % of the climate models, ^[..²³⁵] $T_t^{glob,var}$ can be modeled as an AR(1) process. In the remaining ones either an AR(2) or AR(3) process is chosen. All emulators contain oscillating positive and negative AR coefficients with the first coefficient being positive, but they differ in the ^[..²³⁶] magnitude of the **respective** AR coefficients. The associated innovations vary in their standard deviations by a factor of almost ^[..²³⁷] **three** (0.06–0^[..²³⁸] .15 °C).

In the local ^[..²³⁹] response module (Eq. 6), more than ^[..²⁴⁰] **80 %** of the land grid points warm more quickly than the global mean^[..²⁴¹], i.e., ^[..²⁴²] $\beta_s^{trend} > 1$, in **25 out of 40 emulators** (Fig. 7). ^[..²⁴³] Overall, the spread in the ^[..²⁴⁴] β_s^{trend} terms differs substantially between emulators trained on different climate models. The ^[..²⁴⁵] intercept terms β_s^{int} cluster closely around zero in each ^[..²⁴⁶] emulator. The fraction of outlier grid points deviating ^[..²⁴⁷] > 1 °C ^[..²⁴⁸] from 0, and hinting at sub-optimal local fits, exceeds 1 % in only ^[..²⁴⁹] **one of the emulators**. The vast majority of land grid points are **positively correlated with** $T_t^{glob,var}$, i.e., $\beta_s^{var} > 0$, with the minimum fraction of positive correlations amounting to 82 % ^[..²⁵⁰] of the land grid points.

²¹⁵removed: Additionally, in

²¹⁶removed: calibrated

²¹⁷removed: S1–S9

²³³removed: The deterministic

²³⁴removed: trend $T_t^{glob, det}$ at the end of the 21st century ranges between 3.40 and 6.23

²³⁵removed: the global mean temperature variability

²³⁶removed: respective

²³⁷removed: 3 (0.055

²³⁸removed: .148

²³⁹removed: mean

²⁴⁰removed: 50

²⁴¹removed: temperature

²⁴²removed: $\beta_s^{det} > 1$, in each calibrated emulator

²⁴³removed: In 32 out of 40 calibrated emulators, this is the case in even more than 75 % of the land grid points.

²⁴⁴removed: β_s^{det}

²⁴⁵removed: vast majority of land grid points are positively correlated with global temperature variability $T_t^{glob,var}$, i.e., $\beta_s^{var} > 0$, with the minimum fraction of positive correlations amounting to 82 % of the land grid points. The

²⁴⁶removed: calibrated

²⁴⁷removed: substantially, i.e.,

²⁴⁸removed: ,

²⁴⁹removed: 4 of the calibrated emulators and does not exceed 2

²⁵⁰removed: in any

Emulator calibration parameters

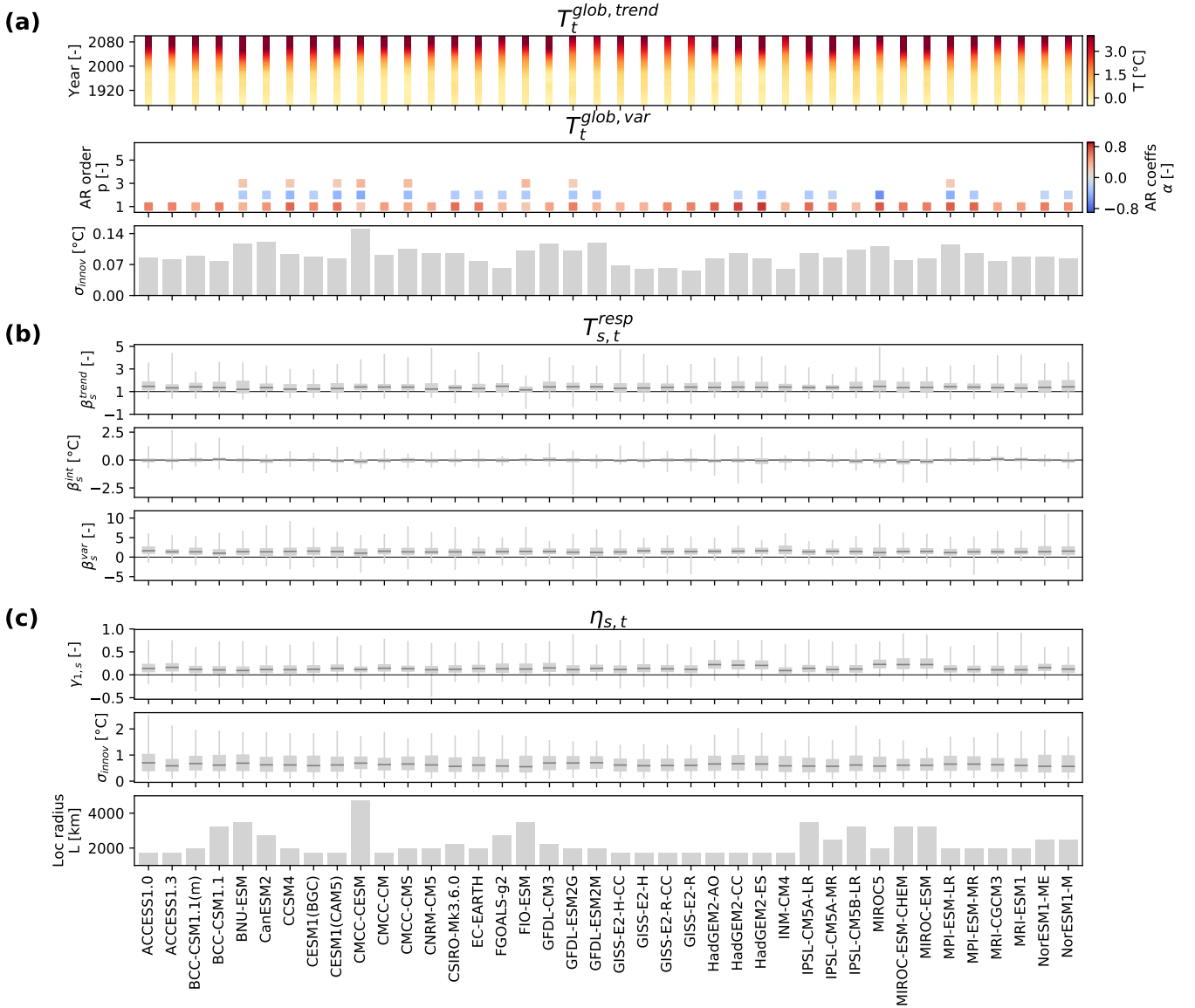


Figure 7. Emulator properties (rows) of the 40 CMIP5 climate models (columns). (a) For the global mean temperature $T_t^{glob, trend}$ and the AR coefficients plus the standard deviation of the innovations of $T_t^{glob, var}$ are depicted. (b) For the local temperature response module, the regression coefficients are shown. (c) For the local residual variability module, the lag-1 AR coefficients, the standard deviations of the innovations, and the localization radii are displayed. Boxplots indicate the median (dark grey line), the interquartile range (grey box), and the full range of values (grey whiskers).

In the local residual variability module (Eq. 7)^[.251], the year-to-year memory contribution $\gamma_{1,s}$ is overall generally small with the 75% quantile lying below 0.25 for 34 out of 40 emulators (Fig. 7).^[.252] Only the six models of the HadGEM and the MIROC family tend to have systematically larger $\gamma_{1,s}$. While the median of the standard deviations of the innovations is similar in all calibrated emulators, the full ranges^[.253] differ substantially, with the maximum^[.254] between 1.3 and 2.5 °C.^[.255] The selected localization radii vary between^[.256] 1750 and 4750 km. Thereby^[.257] 4750 km is a strong outlier with the second highest^[.258] localization radii amounting to 3500 km. Generally, climate models with a coarser native resolution are associated with larger localization radii (not shown).

6.2 Example realizations

^[.259] Figure 8 demonstrates that the emulations nicely capture regional-scale trends and variability in the training and the test runs of the CMIP5 ensemble. The histograms also highlight that the larger sample size of^[.260] the emulations by a factor of 1000 makes it possible to sample the temperature phase space better. The CMIP5^[.261] projections, and thus also the emulations, diverge substantially towards the end of the 21st century in global land and SSA but agree rather well in CEU. At the end of the 21st century, an inter-model spread of roughly 4 °C is observed in global land with models spread out evenly across this space. In SSA, on the other hand, the bulk of the models clusters within a space of 2^[.262] °C and a few outlier models cause the overall CMIP5^[.263] spread of almost 6 °C.

6.3 Quantitative verification

6.3.1 Local trends verification

Correlation between the emulated local trends and the true climate model runs is very high in both training and test runs in all CMIP5^[.264] models, indicating that the forced trends are successfully extracted from each training run (Fig. 9). For

²⁵¹ removed: module, the fraction of grid points with different AR orders varies strongly across the emulators but they agree that considering memory terms of more than one year is only seldom useful

²⁵² removed: In fact, for 36 out of 40 models, either an AR(0) or an AR(1) process is chosen in more than 85 % of all grid points

²⁵³ removed: differs

²⁵⁴ removed: lying between 1.36

²⁵⁵ removed: C and 2.53 °

²⁵⁶ removed: 1500 and 4500

²⁵⁷ removed: 4500

²⁵⁸ removed: selected localization radius amounting to 3000

²⁵⁹ removed: Overall, time series of the emulated quantiles follow the original ones closely but are much smoother due to the

²⁶⁰ removed: 8000 realizations as opposed to 40 runs (Fig. 8). However, in some time instances, the actual

²⁶¹ removed: runs diverge more from the emulated runs indicative of physical processes not accounted for in the emulator set-up. The most notable example occurs approximately between 1960 and 1990 where the median and especially the 5

²⁶² removed: % quantile of the emulated ensemble are warmer than the

²⁶³ removed: ensemble ones in global land and CEU. This is likely related to a response to tropospheric aerosol forcing which is not accounted in our framework (see discussion in Sect. 7.1). The

²⁶⁴ removed: projections,

Regionally averaged temperature time series

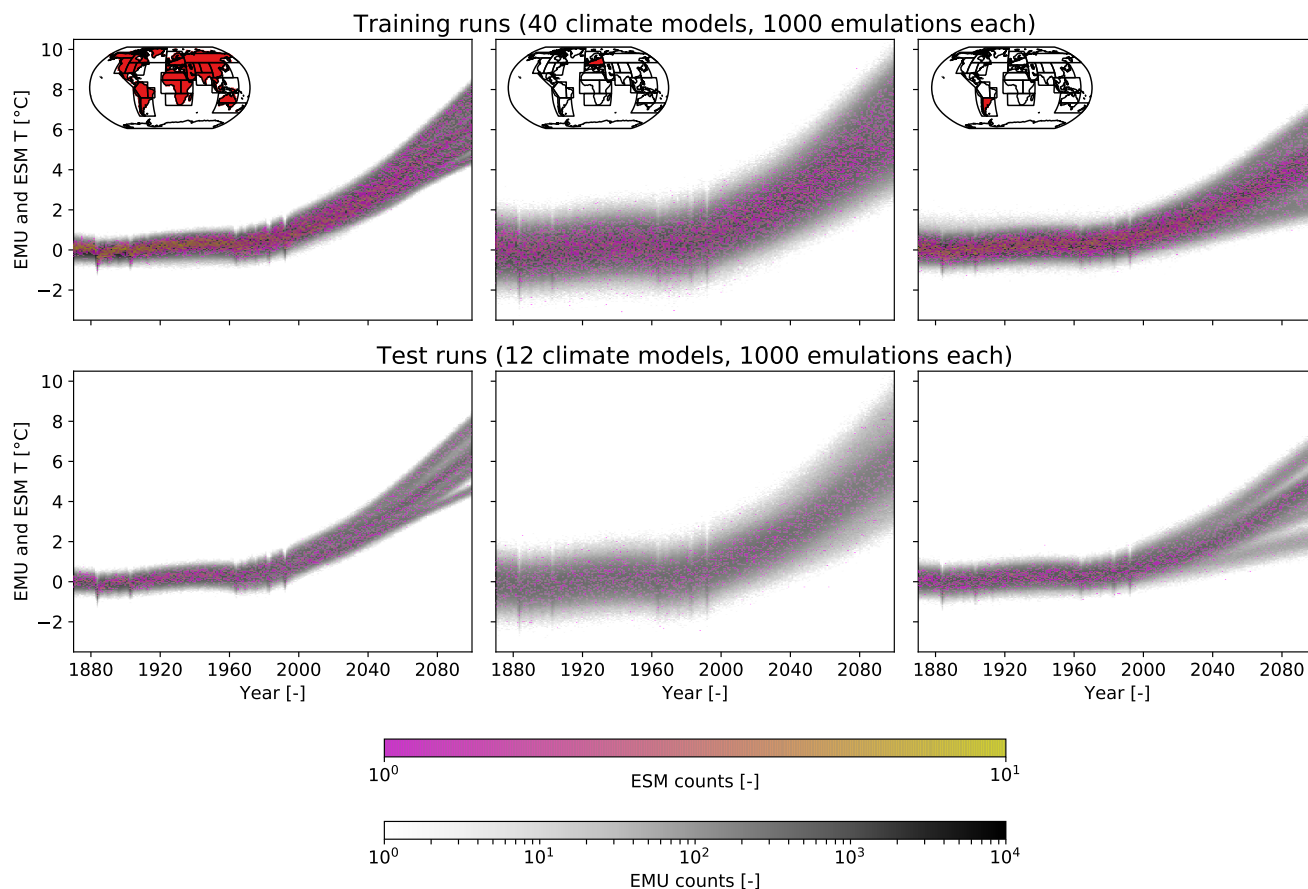


Figure 8. Regionally averaged time series as 2D histograms for 40 CMIP5 model training runs and 1000 emulations per model (top row) and for 12 CMIP5 models with one test run and 1000 emulations per model (bottom row). For the CMIP5 model runs a colormap from pink to yellow is employed and for the emulations a grey-scale is used. The regions are from left to right: global land, CEU, and SSA.

the climate model with test runs, these correlations are nearly identical for each individual test and the training runs. The smallest correlation coefficient is 0.90, the highest one 0.97.

350 6.3.2 Local variability verification

To evaluate the local variability at the grid-point level, lag-1 temporal autocorrelations and standard deviations are considered (Fig. 10). The lag-1 temporal autocorrelation is a rather noisy parameter to estimate and the median correlations between emulations and the training run lie between 0.67 and [\[.265\]](#)0.92. Generally, the correlation of the lag-1 autocorre-

²⁶⁵removed: thus also the emulations, diverge substantially towards the end of the 21st century in global land and SSA

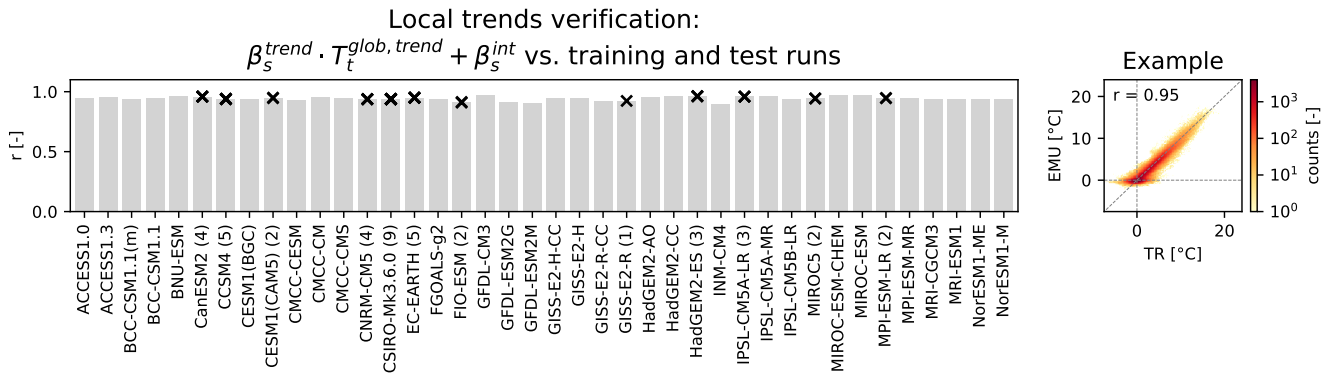


Figure 9. Local trends verification for the CMIP5 models by means of Pearson correlation between the emulated local trends and the training runs (grey bars). The example shows the associated 2D histogram for CESM1(CAM5). For the CMIP5 models with test runs, the correlation between the emulated local trends and each individual test runs is indicated by a black cross. Since these correlations are nearly identical for each test run of a specific climate model, the individual black crosses can visually not be distinguished from one another. For all climate model with test runs, the number of available test runs is given in brackets after the model name.

355 correlations between test and training runs is smaller than the one between emulations and training runs, implying a tendency to overfit this parameter. The correlation between the standard deviations of the emulations and the training run is never below 0.98. The correlation between test and training runs is almost identical to the one between emulations and training runs. Thus, at the grid-point level the emulations reliably reproduce the stochastic variability of climate model runs.

To evaluate the spatial cross-correlations between grid points, three geographical bands are considered (Fig. 11). At all spatial scales, cross-correlations between test and training runs are higher than correlations between emulations and training runs. This is a direct consequence of the regularization which dampens covariances between grid points as a function of distance and is thus inherent to the emulator's design. In a radius of up to 2000 km, the emulators perform best and co-variations between grid points are generally well reproduced. The medians of the correlations between the emulations and the training runs span from 0.85 to 0.98. Plotting an individual example emulation against its associated training run clearly shows the dampening of the cross-correlations in the regularized emulations. Emulations of climate models with larger localization radii (Fig. [..²⁶⁶] 7) have by design a larger correlation with their respective training runs (Fig. 11). In a radius between 2000 and 15000 km, the [..²⁶⁷] emulators perform the least well since there, cross-correlations in the emulations are strongly dampened with the medians of the correlations between emulations and training runs ranging from 0.17 to 0.82. For long-range distances beyond 15000 km, medians lie between 0.20 and 0.93. For all distances beyond 2000 km, there are large inter-model differences in the ability of the emulators to reproduce cross-correlations between grid points. Also correlations between spatial cross-correlations of test and training runs are generally lower

360

365

370

²⁶⁶removed: 8) . In CEU, on the other hand,

²⁶⁷removed: modeled warming rates are more similar. As expected, the difference in smoothness between the emulated quantiles and the CMIP5 ones is most pronounced in regions affected by strong local variability such as CEU

Local variability verification of correlation of grid-point level metrics:

$$\beta_s^{var} \cdot T_t^{glob,var} + \eta_{s,t} \text{ and test runs vs. training runs}$$

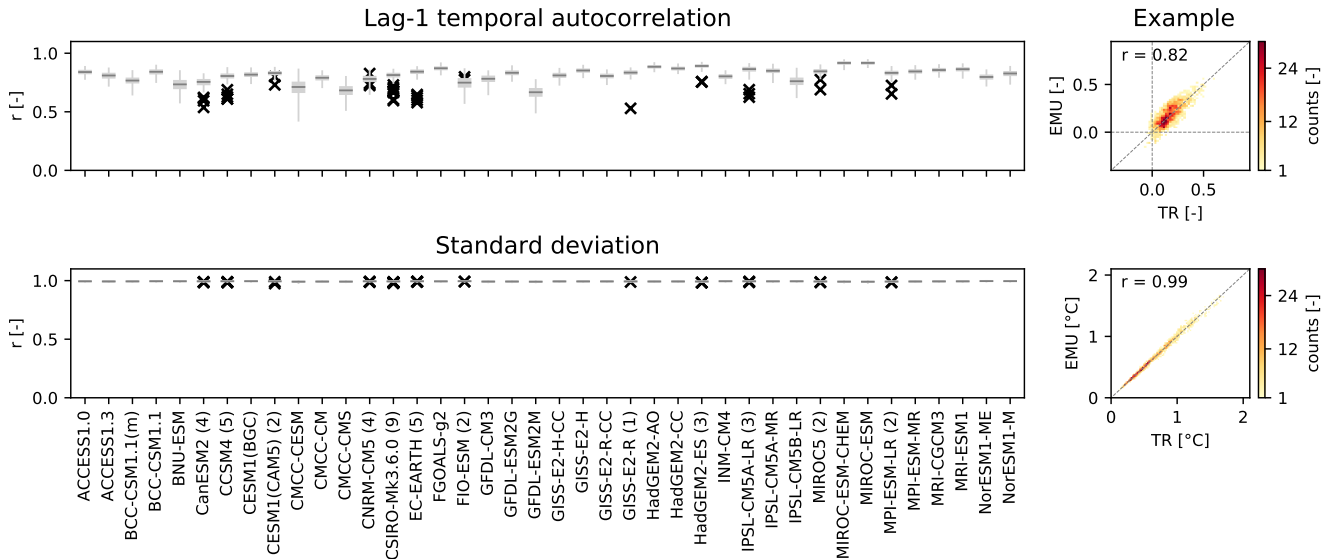


Figure 10. Local variability verification for the CMIP5 models (columns) by means of Pearson correlation of grid-point level lag-1 temporal autocorrelations (top row) and standard deviations (bottom row) between the 1000 individual emulations and the training runs (boxplots). The examples show the associated 2D histograms for a single emulation and the training run of CESM1(CAM5). For the CMIP5 models with test runs, the correlation between the quantity in the training run and in each individual test run is indicated by a black cross. For all climate model with test runs available, the number of test runs is given in brackets after the model name.

and exhibit more inter-model differences at distances beyond 2000 km highlighting that it is more difficult to estimate far reaching spatial cross-correlation based on single ESM runs. Generally, the emulations perform better and are more comparable to test runs at distances beyond 15000 km than between 2000 and 15000 km, which is likely due to that fact that the global correlation pattern induced by the global mean temperature variability serves a more important driver for the longest-range correlations.

375

6.4 [..²⁸⁵]

6.3.1 [..²⁸⁶]Regional-scale emulation verification

[..²⁸⁷] When considering full emulations, i.e., the local trends plus the local variability, the median is successfully emulated but the emulations are a bit underdispersive compared to the training run for the vast majority of CMIP5 models and SREX

²⁸⁵removed: Quantitative verification

²⁸⁶removed: Approximating CMIP5 – the in-sample

²⁸⁷removed: To quantitatively assess the emulators' ability to capture runs they were calibrated on, in-sample verification on the training runs is carried out. For the vast majority of CMIP5 models and SREX regions

Local variability verification of correlation of cross-correlations between grid points:

$$\beta_s^{var} \cdot T_t^{glob, var} + \eta_{s,t} \text{ and test runs vs. training runs}$$

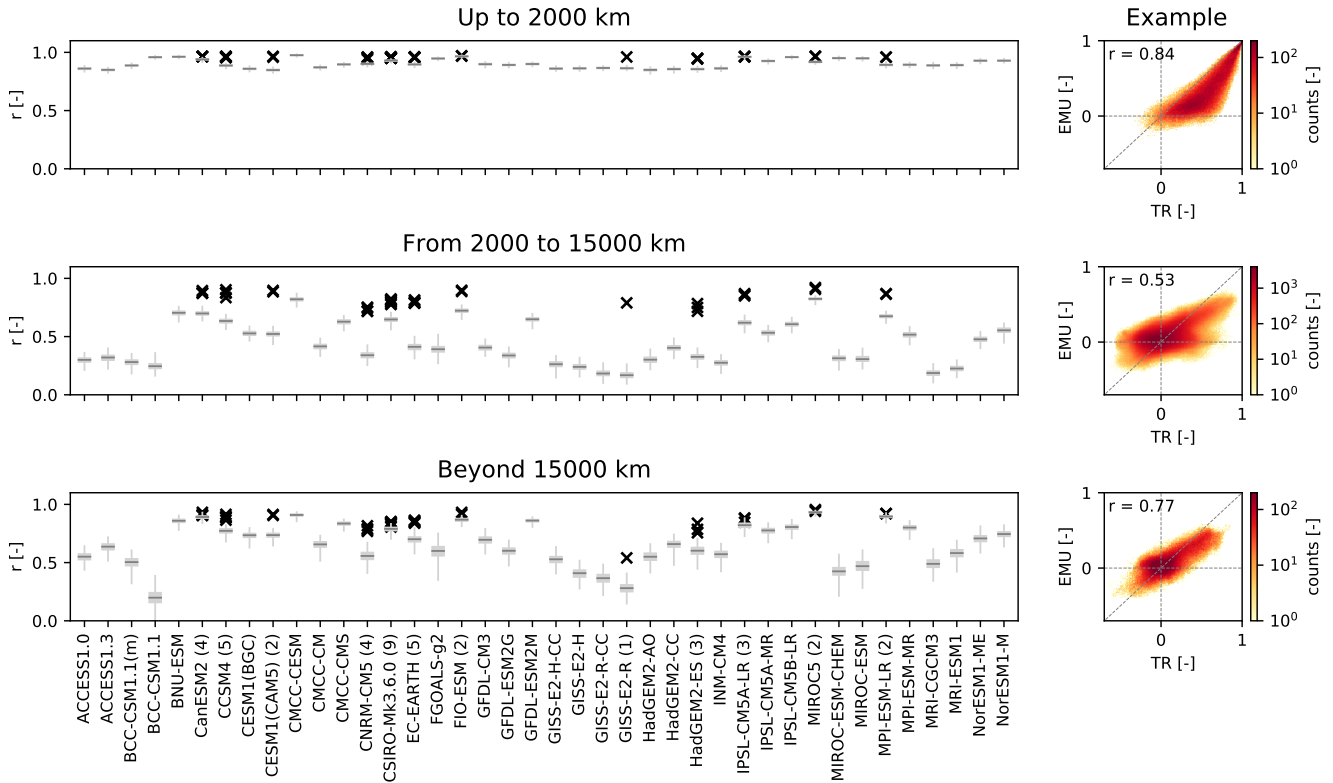


Figure 11. [..²⁶⁸] Local variability verification for the CMIP5 models ([..²⁶⁹] columns) [..²⁷⁰] by means of Pearson correlation of cross-correlations between grid points in three geographical bands (rows) between the 1000 individual emulations and the training runs ([..²⁷¹] boxplots) [..²⁷²]. The geographical bands cover distances below 2000 km, [..²⁷³] between 2000–15000 km, and [..²⁷⁴] beyond 15000 km. The examples show the associated 2D histograms for a single emulsion and the training run of CESM1(CAM5). For the [..²⁷⁵] CMIP5 models with test runs, the [..²⁷⁶] correlation between the [..²⁷⁷] quantity in the [..²⁷⁸] training run and in each individual test run is [..²⁷⁹] indicated by a black cross. For [..²⁸⁰] all climate model with test runs available, the [..²⁸¹] number of test runs is [..²⁸²] given in [..²⁸³] brackets after the [..²⁸⁴] model name.

380 regions (Fig 12). [..²⁸⁸] The emulations tend to be more reliable for climate models [..²⁸⁹] with larger localization radii [..²⁹⁰] (Fig. 7). In North Asia (NAS), [..²⁹¹] the underdispersion is strongest for most models [..²⁹²] (Fig. 12). The only region

²⁸⁸ removed: Generally, the emulations are

²⁸⁹ removed: in which

²⁹⁰ removed: were selected.

²⁹¹ removed: said

²⁹² removed: .

where the emulations are fully reliable is global land. The underdispersion on the SREX regional scales is related to [..²⁹³]the regularization which dampens covariances between grid points as a function of distance between them and is thus inherent to the [..²⁹⁴]

385 [..²⁹⁵]

6.3.2 [..²⁹⁶]

[..²⁹⁷]emulator’s design. The results are qualitatively similar for the test runs but, as expected, the deviations from the emulated quantiles tend to be larger in magnitude than [..²⁹⁸]for the training runs. For most climate models, the strongest deviations in the median of the test runs are observed in global land, Canada/Greenland/Iceland (CGI), and Southern Australia (SAU). [..²⁹⁹]Out of all climate models, the least optimal fit is obtained for MIROC5 [..³⁰⁰]with the emulated median being systematically warmer than the [..³⁰¹]training and especially the test run medians in many regions.

390

7 Discussion

7.1 Emulator design choices and their [..³¹³]implications

[..³¹⁴]

395 7.1.1 Modular framework

A modular framework is chosen for the climate model emulation because of its manifold advantages. First, the calibrated parameters of each emulator module can be used for climate model inter-comparison over a wide range of scales since they can be readily visualized and easily interpreted (Sects. 5.1 and 6.1). Second, the modular framework renders it straightforward to substitute each emulator module with approaches other than the ones chosen here. For example, alternative approaches for the global mean temperature trend [..³¹⁵](e.g., Meinshausen et al., 2011), for the local response

400

²⁹³removed: regularization and

²⁹⁴removed: emulators’ design (see discussion in Sect. 7.1). Out of all climate models, the least optimal fit is obtained for MIROC5 with the emulated median being warmer than the one of the training run in many regions.

²⁹⁵removed: 5 % (left), 50 % (middle), and 95 % (right) quantile for the 40 CMIP5 models (rows) and regions (columns). The deviation of the training run from the emulated quantile is given in color. Blue indicates that the emulated quantile is colder than the quantile of the training run, red means that it is warmer.

²⁹⁶removed: Approximating initial-condition ensembles – the out-of-sample verification

²⁹⁷removed: A more challenging task than capturing the training run is mimicking an initial-condition ensemble, the performance of which can be investigated by carrying out an out-of-sample verification with initial-condition ensemble members not seen during training. Qualitatively similar to the the in-sample verification (Sect. ??), the median is generally well captured but the emulations are a bit underdispersive in regional averages (Fig. ??). However

²⁹⁸removed: in the in-sample verification

²⁹⁹removed: Only

³⁰⁰removed: , the median of the ensemble of emulations is substantially and

³⁰¹removed: climate model runs

³¹³removed: effects

³¹⁴removed: Here, the deterministic

³¹⁵removed: $T_t^{glob,det}$ (Eq. 3) is approximated by

Regional-scale verification: deviation of climate model runs from emulated quantiles

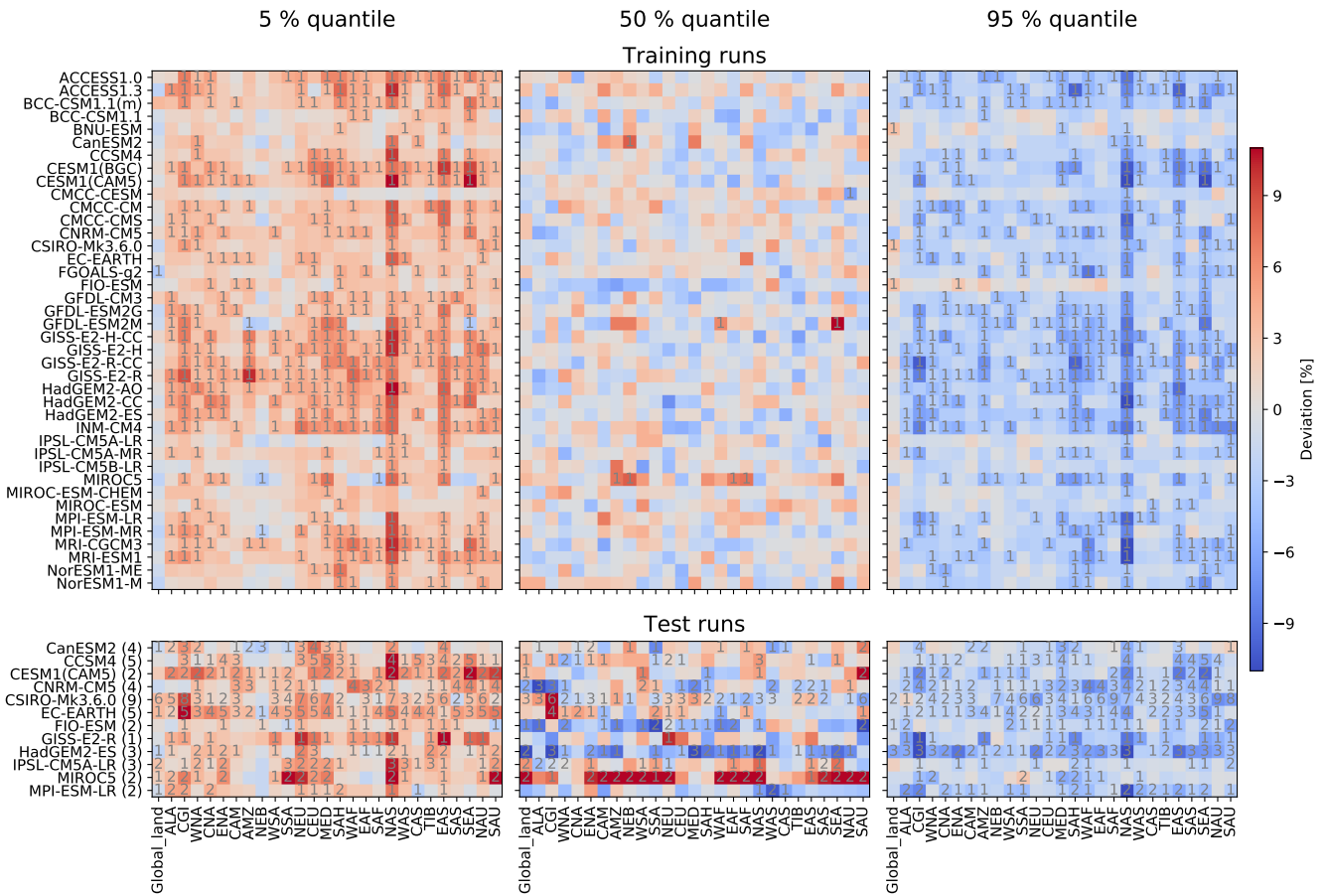


Figure 12. Deviation of climate model runs from the emulated 5 % (left), 50 % (middle), and 95 % (right) quantile for [..³⁰²]CMIP5 models (rows) [..³⁰³]and regions (columns). The [..³⁰⁴]emulated quantile is computed based on 1000 emulations per climate model. The [..³⁰⁵]deviation of [..³⁰⁶]the climate model [..³⁰⁷]run from the emulated quantile is given in color. [..³⁰⁸]Red means that the emulated quantile is [..³⁰⁹]warmer than the quantile of the climate model [..³¹⁰]run, [..³¹¹]blue means that it is [..³¹²]colder. The grey numbers indicate how many climate model run deviations lie outside of the 95 % interval spanned by the deviations of single emulations from the emulated quantiles. If the climate model run lies outside this interval, it is no longer considered indistinguishable from the emulations. The deviation from the training run is shown in the top panel, the average deviation across all available test runs is shown in the bottom panel. The number of test runs averaged across is indicated in brackets behind the model name.

module (e.g., Tebaldi and Arblaster, 2014; Alexeeff et al., 2018), or for the local residual temperature variability (e.g., Link et al., 2019) could be employed. Third, if the modeling task were to change, additional predictors could easily be integrated. For example, precipitation emulation would likely require human-induced aerosol emissions as an additional predictor in the local temperature response module (Frieler et al., 2012).

405 7.1.2 Emulating temperature trends

In this study, an estimate of $T_t^{glob,trend}$ is retrieved with a simple statistical model [..³¹⁶] from the training run (Sect. 4.1.1). However, it could alternatively be considered to obtain $T_t^{glob,trend}$ from a simple energy-balance model (Meinshausen et al., 2011). This would open avenues towards emulating initial-condition [..³¹⁷] ensembles across different $T_t^{glob,trend}$ trajectories and thus different emission scenario pathways.

410 To [..³¹⁸] translate $T_t^{glob,trend}$ into a local [..³¹⁹] temperature in the local response module, a linear approach is chosen ([..³²⁰] Sect. 4.1.2). The thereby obtained regression coefficients β_s^{trend} represent well-known climate phenomena. The enhanced warming over land compared to the global mean (Sutton et al., 2007; Hartmann et al., 2013) at many grid points is captured by $\beta_s^{trend} > 1$ (Figs. 3 and 7). The Arctic amplification (Serreze and Barry, 2011) manifests itself in the large β_s^{trend} values in northern high latitudes (Fig. 3). The overall good performance in capturing local trends is in
415 line with the pronounced linear scaling of regional land temperatures with global mean temperature (Seneviratne et al., 2016; Wartenburger et al., 2017) [..³²¹] and the widely used linear pattern scaling approaches (Mitchell, 2003; Tebaldi and Arblaster, 2014; Lynch et al., 2017; Osborn et al., 2018). [..³²²]

7.1.3 Emulating temperature variability

Spatially coherent local variability is introduced in two emulator modules, namely in the local response module as the
420 local response to $T_t^{glob,var}$ (Sect. [..³²³]
[..³²⁴] 4.1.2) and in the local residual variability module ([..³²⁵] Sect. 4.1.3). The local variability is an essential ingredient in mimicking initial-condition ensembles as visualized by comparing regionally averaged time series of our emulations

³¹⁶removed: accounting for smooth forcing with LOWESS smoothing and for abrupt volcanic forcing with a linear regression to stratospheric optical depth, with global variability $T_t^{glob,var}$ rendering the extraction of $T_t^{glob,det}$ challenging. It is assumed that $T_t^{glob,det}$ is shared by all

³¹⁷removed: ensemble members which is likely not entirely true

³¹⁸removed: downscale the global information

³¹⁹removed: mean response

³²⁰removed: Eq. 6). More complex non-linear methods, specifically, neural networks, were tested as well but resulted in less well-behaved residuals.

³²¹removed: as well as with

³²²removed: The novelty of our local mean response module lies in the explicit separation of the global mean temperature response into a response to the deterministic trend $T_t^{glob,det}$ and to the global variability $T_t^{glob,var}$. This separation allows us to better capture the underlying climate phenomena, namely the local warming and the local co-variation with global variability. Because the local mean response is solely conditioned on global temperature, the emulated quantiles are warmer than the CMIP5 model ones between 1960–1990 in certain regions of the world (see

³²³removed: 6.2). During this time, the rapid increase and subsequent decrease in human-induced aerosol emissions affected surface solar radiation – a phenomenon popularly known as global dimming and brightening – and in turn also surface temperature (Wild, 2012). While our framework could readily be extended to contain tropospheric aerosol predictors, we refrain from it here, as the overall impact on the emulation skill would be small.

³²⁴removed: In the local residual temperature

³²⁵removed: Eq. 7), AR processes are fit to account for temporal autocorrelation. If local trends were not successfully removed before fitting the AR processes, they would introduce artificial long-term memory. However, since AR orders larger than AR(1) are chosen only very rarely, this issue is most likely of minor importance here. A more pronounced effect is caused by regularizing the empirical spatial covariance matrix of the innovations of the local residual temperature variability (Eq. 9) which leads to damped co-variations between grid points and thus underdispersion in regional averages (Sects. 6.3). This explains why the large North Asia region displays the strongest underdispersion and why, generally, more pronounced underdispersion occurs in the CMIP5 models with smaller localization radii, i.e., in the more strongly regularized ones. Global land is the only region not affected by underdispersion since

with simple pattern scaling results which contain no local variability module (Fig. S11). In this study, and all other studies cited in the following paragraphs, the local temperature variability is assumed to be stationary in time which is not fulfilled everywhere in the business-as-usual greenhouse gas emission scenario [..³²⁶](see Sect. 2.3 and Olonscheck and Notz, 2017).
425 [..³²⁷] $T_t^{glob,var}$ can be regarded as the globally aggregated signal of all physical modes of variability (Sect. [..³²⁸]

7.2 [..³²⁹]

[..³³⁰]

[..³³¹]4.1.1), with the calibrated emulators accounting for memory of up to three years (Fig. 7). While the linear trans-
430 lation of $T_t^{glob,var}$ to a grid-point level temperature response is purely statistical in nature, physically meaningful patterns nevertheless emerge in the β_s^{var} patterns. For example, for many climate models, β_s^{var} tends to resemble an El Niño Southern Oscillation [..³³²]pattern (Trenberth, 1997) with Amazon, Australia, and Alaska co-varying while the South-eastern USA exhibits the opposite temperature sign (Figs. 3 and S2–10). Qualitatively similar results could alternatively be obtained by stochastically generating time series of major physical modes of variability and translating those to the
435 grid-point level (McKinnon and Deser, 2018).

Local residual variability is modeled as an AR process [..³³³]with spatially correlated innovations (Sect. 4.1.3). While several other authors have employed AR models with spatially correlated innovations too (Castruccio and Stein, 2013; Castruccio and Genton, 2016; Bao et al., 2016), they all chose a parametric approach to model the covariance between grid points. However, in this study, a non-parametric approach is employed which retains regional-scale anisotropy in the
440 underlying data.

its variability is driven by the global temperature variability $T_t^{glob,var}$, not by the local innovations $\nu_{s,t}$ which, by construction, average out across the globe. The local residual temperature variability module is based on the assumption that residual local temperature variability is

³²⁶removed: (see Sect. 2 and Olonscheck and Notz, 2017)

³²⁷removed: In-sample verification on the training runs and out-of-sample verification on additional initial-condition members not used during training indicate that the assumptions listed above are overall justified since the emulated ensembles generally capture the ESM runs well

³²⁸removed: 6.3). Only in 1 climate model, namely MIROC5, the in-sample and especially the out-of-sample verification reveal a systematic warm bias in the median in many regions. Visual inspection of regionally averaged time series (not shown) suggests that this lack-of-fit arises because MIROC5 reacts strongly to human-induced aerosol forcing which we do not account for .

³²⁹removed: Climate phenomena emerging from the emulator parameters

³³⁰removed: On longer time scales, the most pronounced climate phenomena emerging from the emulator’s calibration parameters are the stronger warming over land compared to the global mean (Sutton et al., 2007; Hartmann et al., 2013) and the Arctic amplification (Serreze and Barry, 2011). In the emulator, the enhanced warming rate over land, which is especially large in the Arctic region, is captured in the regression coefficient β_s^{det} of the local mean response model (Eq. 6). Additionally, in line with the Arctic amplification, large variability around the deterministic trend is present in the northern high latitudes in the calibrated emulators.

³³¹removed: On shorter time scales, quasi-periodic climate phenomena such as the

³³²removed: (Trenberth, 1997) manifest themselves as memory in the global temperature variability $T_t^{glob,var}$ and as characteristic warm and cold anomalies in regions around the globe. In the emulator, this memory signal is reproduced by modeling $T_t^{glob,var}$

³³³removed: (Eq. 5). It is then propagated to the local scale with the local mean response module (Eq. 6). Planetary-scale atmospheric waves, such as Rossby waves (Holton and Hakim, 2013), operate on even shorter time scales and on characteristic, physical-based, length scales. We hypothesize that localization radii below 1500 km are not selected by any of the 40 CMIP5 models because they do not allow to reproduce the large-scale temperature responses induced by planetary-scale atmospheric waves.

7.2 The pros and cons of training on single climate model runs

[..³³⁴]We demonstrated that, for [..³³⁵]yearly temperature at grid-point to regional scales, training on a single run per climate model is sufficient to learn key [..³³⁶]properties of the climate system of this climate model. Early results furthermore indicated that also larger single-model initial-condition ensembles, in that case a 21-member CESM ensemble, can be successfully emulated when training on a single ESM run (Beusch et al., 2018). Since a single run was submitted for the majority of [..³³⁷]climate models participating in CMIP5 [..³³⁸]for the emission pathway considered here, requiring only one run to train the emulator [..³³⁹]gives the opportunity to emulate a much larger multi-model ensemble and thus to have the resulting super-ensemble account for more inter-model uncertainty. Nevertheless, [..³⁴⁰]it is not possible to reproduce the characteristics of a true ESM at all spatial and temporal scales when training on a single run. To obtain the best possible [..³⁴¹]emulations to be used e.g., for uncertainty propagation in climate impact or integrated assessment models, it is thus advisable to employ all available runs for training instead of just a single one for each climate model. [..³⁴²]When training on multiple runs, the parameters of the emulator can be estimated more robustly, which, among other things, results in a larger localization radius and thus the ability to reproduce farther reaching spatial cross-correlations between grid points.

7.3 [..³⁴³]

455 [..³⁴⁴]

7.3 Large single-model initial-condition vs. large multi-model ensembles

Our results highlight fundamental differences between large single-model initial-condition ensembles (Deser et al., 2012; Fischer et al., 2013; Kay et al., 2015; Leduc et al., 2019; Maher et al., 2019) and large multi-model ensembles (Meehl et al., 2007;

³³⁴removed: As discussed above, we demonstrate

³³⁵removed: temperature at the spatio-temporal scale considered here

³³⁶removed: underlying statistical

³³⁷removed: modeling groups

³³⁸removed: submitted a single run

³³⁹removed: on

³⁴⁰removed: in order to

³⁴¹removed: emulators

³⁴²removed: Given the assumption that an initial-condition ensemble shares the statistical properties estimated during the emulator calibration, the more realizations are considered during training, the more robustly underlying trends as well as co-variations between grid points can be estimated.

³⁴³removed: The advantages of a modular framework

³⁴⁴removed: The advantages of the modular framework approach chosen in this study are manifold. First, the calibrated parameters of each emulator module serve as an interesting approach for climate model inter-comparison over a wide range of scales which can be clearly visualized and easily interpreted (Sects. 5.1 and 6.1). Second, the modular framework renders it straightforward to substitute each emulator module with approaches other than the ones chosen here. For example, with alternative approaches for the deterministic global mean temperature trend (e.g., Meinshausen et al., 2011), for the local mean temperature model (e.g., Tebaldi and Arblaster, 2014; Alexeeff et al., 2018), or for the local residual temperature variability (e.g., Link et al., 2019). Third, if the modeling task were to change, additional predictors could easily be integrated. For example, precipitation emulation would likely require human-induced aerosol emissions as an additional predictor in the local mean response model (Frieler et al., 2012).

Taylor et al., 2012; Eyring et al., 2016). While multi-model ensembles are imperfect, with several ESMs [..³⁴⁵] exhibiting dependencies (Knutti, 2010; Bishop and Abramowitz, 2013; Sanderson et al., 2015; Abramowitz et al., 2019), multi-model uncertainty nevertheless clearly exceeds single-model initial-condition uncertainty at the yearly scale for temperature (Sect. 5.3). ESMs contained within CMIP5 differ [..³⁴⁶]substantially across a broad range of scales [..³⁴⁷]and thus sample different phase spaces in projections [..³⁴⁸]which renders it necessary to train an emulator on each climate model to approximate the CMIP5 ensemble. A single-model initial-condition ensemble, on the other hand, can be successfully mimicked on grid-point to regional scales by training on a single [..³⁴⁹]ESM run (Sects. 5 and [..³⁵⁰]6). While this lies beyond the scope of this study, the developed emulator could additionally serve as a novel tool to address the challenge of inter-model dependencies. Differences between climate models could be quantified in terms of their [..³⁵¹]emulator parameters and subsequently, a subset of models with sufficiently divergent parameters could be selected to base projections on. Additionally, observations could be used to constrain the emulated ensemble by providing validation measures for the emulator parameters.

470 8 Conclusions and outlook

We [..³⁵²]introduce a modular framework for climate model emulation of yearly land temperatures and [..³⁵³]present a specific, computationally cheap implementation, which [..³⁵⁴]can create plausible temperature field time series within seconds based on a single climate model training run. Our emulator consists of (i) a global mean temperature module, (ii) a local [..³⁵⁵]temperature response module, and (iii) a local residual temperature variability module. The global mean temperature module contains a [..³⁵⁶]global mean temperature trend which is shared by all emulations and a global mean temperature variability term which is modeled as an AR process and varies between individual emulations. The local [..³⁵⁷]response module is linear in nature and consists of a separate response to the [..³⁵⁸]global mean temperature trend and the global mean temperature variability. The local residual variability module generates spatio-temporally correlated fields by means of locally fit AR(1) processes with spatially correlated innovations.

480 [..³⁵⁹]Since emulators approximate complex ESMs in a simplified manner, they are not able to accurately reproduce all spatio-temporal ESM characteristics. The emulator presented here, e.g., dampens co-variations between grid points as

³⁴⁵removed: showing

³⁴⁶removed: sufficiently

³⁴⁷removed: to

³⁴⁸removed: ,

³⁴⁹removed: available

³⁵⁰removed: ??

³⁵¹removed: calibrated

³⁵²removed: introduced

³⁵³removed: presented

³⁵⁴removed: makes it possible to

³⁵⁵removed: mean temperature

³⁵⁶removed: deterministic trend as well as a stochastic

³⁵⁷removed: mean temperature model

³⁵⁸removed: the deterministic global

³⁵⁹removed: With emulators trained on single runs of four selected ESMs, we showed that the emulations

a function of distance in the local residual variability module due to regularization. Thus, our emulator reliably reproduces climate model variability at the grid-point level, but the emulations are increasingly underdispersive for larger regional averages and intermediate-range spatial teleconnections cannot be accounted for. This caveat could be addressed by further improving the local residual variability module implementation with a focus on such teleconnections. Alternatively, training on several ESM runs would increase the robustness of the estimated parameters and make it possible to reproduce farther-reaching teleconnections within the current emulator setup. Nevertheless, calibrating our emulator on a single training run is sufficient to generate emulations which are visually indistinguishable from ³⁶⁰true ESM runs.

Inherent inter-ESM differences in ³⁶¹warming trends and spatio-temporal variability make it necessary to calibrate a separate emulator for each ³⁶²one of the 40 considered CMIP5 models³⁶³. The resulting emulations successfully approximate the training run for each climate model on grid-point to regional scales. For CMIP5 models with more than one initial-condition ensemble member, it was furthermore demonstrated that ³⁶⁴the ensemble of emulations is generally able to mimic true climate model initial-condition ensembles ³⁶⁵at these scales. Hence, we argue that to sample climate signal uncertainty for yearly temperature at grid-point to regional scales, it is more advantageous to invest computational resources into generating multi-model ensembles rather than large single-model ensembles, since the latter can be readily approximated by our emulator.

Super-ensembles such as the one generated in this study, which contains 1000 emulations per climate model, are expected to be ³⁶⁶particularly helpful in regions with large interannual variability. There, the very ³⁶⁷sparse sampling of the temperature phase space by the CMIP5 ensemble ³⁶⁸may result in biased conclusions when solely employing the CMIP5 ensemble as an input to impact or integrated assessment models which estimate the ³⁶⁹effect of climate signal uncertainty on their quantity of interest.

³⁷⁰The emulator is designed to be flexible enough to emulate whatever climate model run it is provided with. Hence, it is not part of the emulator's tasks to judge the realism of individual climate models. Instead, the choice of considered

³⁶⁰removed: ESM runs not employed during training and that inherent

³⁶¹removed: deterministic

³⁶²removed: considered climate model. Thus, we proceeded to emulate

³⁶³removed: by calibrating an emulator on a single run from each climate model. The ensembles of emulations closely resemble the training runs

³⁶⁴removed: , generally, the ensembles of emulations are

³⁶⁵removed: . However, a trade-off between parameter estimation robustness and ensemble reliability results in emulations which are, by construction, slightly underdispersive in regional averages. Hence, the generated ensembles overall provide a conservative estimate of regional-scale variability around the mean trend. While the super-ensemble of this study contains 8000 runs – 200 per CMIP5 climate model– it could readily be extended to an arbitrary larger number due to the negligible computational cost of the emulator. Such super-ensembles

³⁶⁶removed: especially

³⁶⁷removed: noisy quantiles of the

³⁶⁸removed: could result in misleading

³⁶⁹removed: impact

³⁷⁰removed: Hence, we argue that to sample climate signal uncertainty for yearly temperature at the spatial scale considered in this study, it is more advantageous to invest computational resources in generating multi-model ensembles rather than large single-model ensembles, since the latter can be readily mimicked by our emulator based on a single ESM run . However, when the goal is to obtain the best possible emulator, we nevertheless advise to train on all available runs with differing initial conditions to enhance the robustness of the estimated statistical parameters. In addition, it could be considered to combine findings

ESMs will depend on the scope of different applications. For example, results from emergent constraints analyses (e.g.,
505 Hall and Qu, 2006; Eyring et al., 2019) could be combined with the implementation of [..³⁷¹] an emulator to derive a super-ensemble based on an observationally-constrained set of ESMs.

[..³⁷²] On the other hand, the emulator parameters could themselves be used as potential constraints that can also be derived from observations. Additionally, the emulator parameters can be regarded as an ESM-specific "model ID" [..³⁷³] which provides an interesting avenue for climate model inter-comparison across a wide range of scales. Inter-model
510 differences can be [..³⁷⁴] readily visualized for every emulator module resulting in comprehensible scale-dependent insights into the underlying [..³⁷⁵] properties of each climate model.

Future work could focus on [..³⁷⁶] extending the emulator to simultaneously generate [..³⁷⁷] multivariate output. Furthermore, it would be interesting to investigate how transferable an emulator trained on a specific greenhouse gas emission scenario is to other emission pathways and which modules would need to be modified to account for inter-scenario differences.

515 In conclusion, in this study we have presented a novel ESM emulator that can be trained to represent separate ESMs based on single realizations of the respective ESMs, and which has been shown to be able to emulate and expand multi-model ensembles such as CMIP5. We expect that the developed emulator can serve as training ground for investigating the phase space of multi-model ensembles in new applications, e.g. related to the derivation of emissions scenarios or the assessment of impacts under different emissions pathways.

520 *Data availability.* The employed CMIP5 data are available from the public CMIP archive at <https://esgf-node.llnl.gov/projects/esgf-llnl/>. The stratospheric aerosol optical depth data are provided by NASA and available at <https://data.giss.nasa.gov/modelforce/strataer/>.

Appendix A

Author contributions. LB, LG, and SIS designed the study, based on an initial idea from SIS. LB carried out the analysis and drafted the text. LG provided statistical support for the analysis. All authors contributed to interpreting the results and refining the text.

525 *Competing interests.* The authors declare that they have no conflict of interest.

³⁷¹removed: a climate

³⁷²removed: The calibrated emulator parameters can be

³⁷³removed: and provide

³⁷⁴removed: directly

³⁷⁵removed: statistical

³⁷⁶removed: introducing additional physical aspects into our statistical emulator by e.g., explicitly modeling modes of interannual variability such as the El Niño Southern Oscillation and the Arctic Oscillation. Furthermore, the emulator could be extended

³⁷⁷removed: temperature and precipitation fields. Lastly

Table A1. List of the 40 employed CMIP5 models, the modeling groups providing them, and the number of initial-condition ensemble members used.

Model	Modeling Center (or Group)	Runs
ACCESS1.0	Commonwealth Scientific and Industrial Research Organization (CSIRO) and Bureau of Meteorology (BOM), Australia	1
ACCESS1.3	Commonwealth Scientific and Industrial Research Organization (CSIRO) and Bureau of Meteorology (BOM), Australia	1
BCC-CSM1.1(m)	Beijing Climate Center, China Meteorological Administration	1
BCC-CSM1.1	Beijing Climate Center, China Meteorological Administration	1
BNU-ESM	College of Global Change and Earth System Science, Beijing Normal University	1
CanESM2	Canadian Centre for Climate Modeling and Analysis	5
CCSM4	National Center for Atmospheric Research	6
CESM1(BGC)	Community Earth System Model Contributors	1
CESM1(CAM5)	Community Earth System Model Contributors	3
CMCC-CESM	Centro Euro-Mediterraneo per I Cambiamenti Climatici	1
CMCC-CM	Centro Euro-Mediterraneo per I Cambiamenti Climatici	1
CMCC-CMS	Centro Euro-Mediterraneo per I Cambiamenti Climatici	1
CNRM-CM5	Centre National de Recherches Météorologiques / Centre Européen de Recherche et Formation Avancée en Calcul Scientifique	5
CSIRO-Mk3.6.0	Commonwealth Scientific and Industrial Research Organization in collaboration with Queensland Climate Change Centre of Excellence	10
EC-EARTH	EC-EARTH consortium	6
FGOALS-g2	LASG, Institute of Atmospheric Physics, Chinese Academy of Sciences and CESS, Tsinghua University	1
FIO-ESM	The First Institute of Oceanography, SOA, China	3
GFDL-CM3	NOAA Geophysical Fluid Dynamics Laboratory	1
GFDL-ESM2G	NOAA Geophysical Fluid Dynamics Laboratory	1
GFDL-ESM2M	NOAA Geophysical Fluid Dynamics Laboratory	1
GISS-E2-H-CC	NASA Goddard Institute for Space Studies	1
GISS-E2-H	NASA Goddard Institute for Space Studies	1
GISS-E2-R-CC	NASA Goddard Institute for Space Studies	1
GISS-E2-R	NASA Goddard Institute for Space Studies	2
HadGEM2-AO	National Institute of Meteorological Research/Korea Meteorological Administration	1
HadGEM2-CC	Met Office Hadley Centre	1
HadGEM2-ES	Met Office Hadley Centre (additional realizations contributed by Instituto Nacional de Pesquisas Espaciais)	4
INM-CM4	Institute for Numerical Mathematics	1
IPSL-CM5A-LR	Institut Pierre-Simon Laplace	4
IPSL-CM5A-MR	Institut Pierre-Simon Laplace	1
IPSL-CM5B-LR	Institut Pierre-Simon Laplace	1
MIROC5	Atmosphere and Ocean Research Institute (The University of Tokyo), National Institute for Environmental Studies, and Japan Agency for Marine-Earth Science and Technology	3
MIROC-ESM-CHEM	Japan Agency for Marine-Earth Science and Technology, Atmosphere and Ocean Research Institute (The University of Tokyo), and National Institute for Environmental Studies	1
MIROC-ESM	Japan Agency for Marine-Earth Science and Technology, Atmosphere and Ocean Research Institute (The University of Tokyo), and National Institute for Environmental Studies	1
MPI-ESM-LR	Max-Planck-Institut für Meteorologie (Max Planck Institute for Meteorology)	3
MPI-ESM-MR	Max-Planck-Institut für Meteorologie (Max Planck Institute for Meteorology)	1
MRI-CGCM3	Meteorological Research Institute	1
MRI-ESM1	Meteorological Research Institute	1
NorESM1-ME	Norwegian Climate Centre	1
NorESM1-M	Norwegian Climate Centre	1

Acknowledgements. We acknowledge partial support from the H2020 CRESCENDO project (grant agreement 641816) and the ERC DROUGHT-HEAT project (grant agreement 617518). We furthermore acknowledge the World Climate Research Program’s Working Group on Coupled

Modeling, which is responsible for the Coupled Model Intercomparison Project (CMIP), and we thank the climate modeling groups (listed in Table A1 of this paper) for producing and making available their model output. We are additionally indebted to Urs Beyerle and Jan
530 Sedláček for retrieving and pre-processing the CMIP5 data. [³⁷⁸]Moreover, we would like to thank Julien Brajard, Loris Foresti, and Vincent Humphrey for their valuable input on different modules of our emulator. Lastly, we thank Robert Link and the two anonymous reviewers for their useful feedback which helped to improve this study.

³⁷⁸removed: Lastly

References

- Abramowitz, G., Herger, N., Gutmann, E., Hammerling, D., Knutti, R., Leduc, M., Lorenz, R., Pincus, R., and Schmidt, G. A.: ESD Reviews: Model dependence in multi-model climate ensembles: Weighting, sub-selection and out-of-sample testing, *Earth System Dynamics*, 10, 91–105, <https://doi.org/10.5194/esd-10-91-2019>, 2019.
- Alexeeff, S. E., Nychka, D., Sain, S. R., and Tebaldi, C.: Emulating mean patterns and variability of temperature across and within scenarios in anthropogenic climate change experiments, *Climatic Change*, 146, 319–333, <https://doi.org/10.1007/s10584-016-1809-8>, 2018.
- Allen, M. R. and Stott, P. A.: Estimating signal amplitudes in optimal fingerprinting, part I: Theory, *Climate Dynamics*, 21, 477–491, <https://doi.org/10.1007/s00382-003-0313-9>, 2003.
- Bao, J., McInerney, D. J., and Stein, M. L.: A spatial-dependent model for climate emulation, *Environmetrics*, 27, 396–408, <https://doi.org/10.1002/env.2412>, 2016.
- Beusch, L., Gudmundsson, L., and Seneviratne, S. I.: Emulating Earth System Model Temperatures, in: Proceedings of the 8th International Workshop on Climate Informatics: CI 2018, edited by Chen, C., Cooley, D., Runge, J., and Szekely, E., pp. 41–44, NCAR Technical Note, Boulder, <https://doi.org/10.5065/D6BZ64XQ>, 2018.
- Bishop, C. H. and Abramowitz, G.: Climate model dependence and the replicate Earth paradigm, *Climate Dynamics*, 41, 885–900, <https://doi.org/10.1007/s00382-012-1610-y>, 2013.
- Brown, P. T., Li, W., Cordero, E. C., and Mauget, S. A.: Comparing the model-simulated global warming signal to observations using empirical estimates of unforced noise, *Scientific Reports*, 5, 1–9, <https://doi.org/10.1038/srep09957>, 2015.
- Carrassi, A., Bocquet, M., Bertino, L., and Evensen, G.: Data assimilation in the geosciences: An overview of methods, issues, and perspectives, *Wiley Interdisciplinary Reviews: Climate Change*, 9, 1–79, <https://doi.org/10.1002/wcc.535>, 2018.
- Castruccio, S. and Genton, M. G.: Compressing an Ensemble With Statistical Models: An Algorithm for Global 3D Spatio-Temporal Temperature, *Technometrics*, 58, 319–328, <https://doi.org/10.1080/00401706.2015.1027068>, 2016.
- Castruccio, S. and Stein, M. L.: Global space-time models for climate ensembles, *Annals of Applied Statistics*, 7, 1593–1611, <https://doi.org/10.1214/13-AOAS656>, 2013.
- Castruccio, S., McInerney, D. J., Stein, M. L., Crouch, F. L., Jacob, R. L., and Moyer, E. J.: Statistical emulation of climate model projections based on precomputed GCM runs, *Journal of Climate*, 27, 1829–1844, <https://doi.org/10.1175/JCLI-D-13-00099.1>, 2014.
- Cressie, N. and Wikle, C. K.: *Statistics for spatio-temporal data*, John Wiley & Sons, 2011.
- Deser, C., Phillips, A., Bourdette, V., and Teng, H.: Uncertainty in climate change projections: The role of internal variability, *Climate Dynamics*, 38, 527–546, <https://doi.org/10.1007/s00382-010-0977-x>, 2012.
- Eyring, V., Bony, S., Meehl, G. A., Senior, C. A., Stevens, B., Stouffer, R. J., and Taylor, K. E.: Overview of the Coupled Model Intercomparison Project Phase 6 (CMIP6) experimental design and organization, *Geoscientific Model Development*, 9, 1937–1958, <https://doi.org/10.5194/gmd-9-1937-2016>, 2016.
- Eyring, V., Cox, P. M., Flato, G. M., Gleckler, P. J., Abramowitz, G., Caldwell, P., Collins, W. D., Gier, B. K., Hall, A. D., Hoffman, F. M., Hurtt, G. C., Jahn, A., Jones, C. D., Klein, S. A., Krasting, J. P., Kwiatkowski, L., Lorenz, R., Maloney, E., Meehl, G. A., Pendergrass, A. G., Pincus, R., Ruane, A. C., Russell, J. L., Sanderson, B. M., Santer, B. D., Sherwood, S. C., Simpson, I. R., Stouffer, R. J., and Williamson, M. S.: Taking climate model evaluation to the next level, *Nature Climate Change*, 9, 102–110, <https://doi.org/10.1038/s41558-018-0355-y>, 2019.

- Fischer, E. M., Beyerle, U., and Knutti, R.: Robust spatially aggregated projections of climate extremes, *Nature Climate Change*, 3, 1033–1038, <https://doi.org/10.1038/nclimate2051>, <http://dx.doi.org/10.1038/nclimate2051>, 2013.
- 570 Frieler, K., Meinshausen, M., Mengel, M., Braun, N., and Hare, W.: A scaling approach to probabilistic assessment of regional climate change, *Journal of Climate*, 25, 3117–3144, <https://doi.org/10.1175/JCLI-D-11-00199.1>, 2012.
- Gaspari, G. and Cohn, S. E.: Construction of correlation functions in two and three dimensions, *Quarterly Journal of the Royal Meteorological Society*, 125, 723–757, <https://doi.org/10.1256/smsqj.55416>, 1999.
- 575 Goodwin, P.: How historic simulation–observation discrepancy affects future warming projections in a very large model ensemble, *Climate Dynamics*, 47, 2219–2233, <https://doi.org/10.1007/s00382-015-2960-z>, 2016.
- Hall, A. and Qu, X.: Using the current seasonal cycle to constrain snow albedo feedback in future climate change, *Geophysical Research Letters*, 33, L03 502, <https://doi.org/10.1029/2005GL025127>, 2006.
- Hartmann, D., Klein Tank, A., Rusticucci, M., Alexander, L., Brönnimann, S., Charabi, Y., Dentener, F., Dlugokencky, E., Easterling, D., 580 Kaplan, A., Soden, B., Thorne, P., Wild, M., and Zhai, P.: Observations: Atmosphere and Surface, in: *Climate Change 2013: The Physical Science Basis. Contribution of Working Group I to the Fifth Assessment Report of the Intergovernmental Panel on Climate Change*, edited by Stocker, T., Qin, D., Plattner, G.-K., Tignor, M., Allen, S., Boschung, J., Nauels, A., Xia, Y., Bex, V., and Midgle, P., chap. 2, pp. 159–254, Cambridge, United Kingdom and New York, NY, USA, <https://doi.org/10.1017/CBO9781107415324.008>, 2013.
- Hawkins, E. and Sutton, R.: The potential to narrow uncertainty in regional climate predictions, *Bulletin of the American Meteorological Society*, 90, 1095–1107, <https://doi.org/10.1175/2009BAMS2607.1>, 2009.
- 585 Herger, N., Sanderson, B. M., and Knutti, R.: Improved pattern scaling approaches for the use in climate impact studies, *Geophysical Research Letters*, 42, 3486–3494, <https://doi.org/10.1002/2015GL063569>, 2015.
- Holden, P. B. and Edwards, N. R.: Dimensionally reduced emulation of an AOGCM for application to integrated assessment modelling, *Geophysical Research Letters*, 37, 1–5, <https://doi.org/10.1029/2010GL045137>, 2010.
- 590 Holden, P. B., Edwards, N. R., Garthwaite, P. H., Fraedrich, K., Lunkeit, F., Kirk, E., Labriet, M., Kanudia, A., and Babonneau, F.: PLASIM-ENTSem v1.0: a spatio-temporal emulator of future climate change for impacts assessment, *Geoscientific Model Development*, 7, 433–451, <https://doi.org/10.5194/gmd-7-433-2014>, 2014.
- Holton, J. R. and Hakim, G. J.: *An introduction to dynamic meteorology*, Elsevier, 5 edn., <https://doi.org/10.1016/C2009-0-63394-8>, 2013.
- Humphrey, V. and Gudmundsson, L.: GRACE-REC: a reconstruction of climate-driven water storage changes over the last century, *Earth System Science Data*, 11, 1153–1170, <https://doi.org/10.5194/essd-11-1153-2019>, 2019.
- 595 Kay, J. E., Deser, C., Phillips, A., Mai, A., Hannay, C., Strand, G., Arblaster, J. M., Bates, S. C., Danabasoglu, G., Edwards, J., Holland, M., Kushner, P., Lamarque, J. F., Lawrence, D., Lindsay, K., Middleton, A., Munoz, E., Neale, R., Oleson, K., Polvani, L., and Vertenstein, M.: The community earth system model (CESM) large ensemble project : A community resource for studying climate change in the presence of internal climate variability, *Bulletin of the American Meteorological Society*, 96, 1333–1349, <https://doi.org/10.1175/BAMS-D-13-00255.1>, 2015.
- 600 King, A. D., Karoly, D. J., and Henley, B. J.: Australian climate extremes at 1.5 °c and 2 °c of global warming, *Nature Climate Change*, 7, 412–416, <https://doi.org/10.1038/nclimate3296>, 2017.
- Knutti, R.: The end of model democracy?, *Climatic Change*, 102, 395–404, <https://doi.org/10.1007/s10584-010-9800-2>, 2010.
- 605 Knutti, R., Masson, D., and Gettelman, A.: Climate model genealogy: Generation CMIP5 and how we got there, *Geophysical Research Letters*, 40, 1194–1199, <https://doi.org/10.1002/grl.50256>, 2013.

- Leduc, M., Mailhot, A., Frigon, A., Martel, J.-L., Ludwig, R., Brietzke, G. B., Giguère, M., Brissette, F., Turcotte, R., Braun, M., and Scinocca, J.: ClimEx project: a 50-member ensemble of climate change projections at 12-km resolution over Europe and north-eastern North America with the Canadian Regional Climate Model (CRCM5), *Journal of Applied Meteorology and Climatology*, 0, <https://doi.org/10.1175/jamc-d-18-0021.1>, 2019.
- 610 Levy, H., Horowitz, L. W., Schwarzkopf, M. D., Ming, Y., Golaz, J.-C., Naik, V., and Ramaswamy, V.: The roles of aerosol direct and indirect effects in past and future climate change, *Journal of Geophysical Research: Atmospheres*, 118, 4521–4532, <https://doi.org/10.1002/jgrd.50192>, 2013.
- Link, R., Snyder, A., Lynch, C., Hartin, C., Kravitz, B., and Bond-Lamberty, B.: Fldgen v1.0: an emulator with internal variability and space – time correlation for Earth system models, *Geoscientific Model Development*, 12, 1477–1489, [https://doi.org/10.5194/gmd-12-](https://doi.org/10.5194/gmd-12-1477-2019)
- 615 1477-2019, 2019.
- Lopez, A., Suckling, E. B., and Smith, L. A.: Robustness of pattern scaled climate change scenarios for adaptation decision support, *Climatic Change*, 122, 555–566, <https://doi.org/10.1007/s10584-013-1022-y>, 2014.
- Lynch, C., Hartin, C., Bond-Lamberty, B., and Kravitz, B.: An open-access CMIP5 pattern library for temperature and precipitation: Description and methodology, *Earth System Science Data*, 9, 281–292, <https://doi.org/10.5194/essd-9-281-2017>, 2017.
- 620 Maher, N., Milinski, S., Suarez-Gutierrez, L., Botzet, M., Kornbluh, L., Takano, Y., Kröger, J., Ghosh, R., Hedemann, C., Li, C., Li, H., Manzini, E., Notz, D., Putrasahan, D., Boysen, L., Claussen, M., Ilyina, T., Olonscheck, D., Raddatz, T., Stevens, B., and Marotzke, J.: The Max Planck Institute Grand Ensemble – Enabling the Exploration of Climate System Variability, *Journal of Advances in Modeling Earth Systems*, 11, 1–21, <https://doi.org/10.1029/2019MS001639>, 2019.
- May, W.: Assessing the strength of regional changes in near-surface climate associated with a global warming of 2°C, *Climatic Change*, 110, 619–644, <https://doi.org/10.1007/s10584-011-0076-y>, 2012.
- 625 McKinnon, K. A. and Deser, C.: Internal variability and regional climate trends in an observational large ensemble, *Journal of Climate*, 31, 6783–6802, <https://doi.org/10.1175/JCLI-D-17-0901.1>, 2018.
- McKinnon, K. A., Poppick, A., Dunn-Sigouin, E., and Deser, C.: An "observational large ensemble" to compare observed and modeled temperature trend uncertainty due to internal variability, *Journal of Climate*, 30, 7585–7598, <https://doi.org/10.1175/JCLI-D-16-0905.1>,
- 630 2017.
- Meehl, G. A., Covey, C., Delworth, T., Latif, M., McAvaney, B., Mitchell, J. F. B., Stouffer, R. J., and Taylor, K. E.: The WCRP CMIP3 multi-model dataset: a new era in climate change research *Bulletin of the American Meteorological Society*, *Bulletin of the American Meteorological Society*, 88, 1383–1394, <https://doi.org/10.1175/BAMS-88-9-1383>, 2007.
- Meinshausen, M., Raper, S. C. B., and Wigley, T. M. L.: Emulating coupled atmosphere-ocean and carbon cycle models with a simpler model, *MAGICC6 - Part 1: Model description and calibration*, *Atmospheric Chemistry and Physics*, 11, 1417–1456, [https://doi.org/10.5194/acp-](https://doi.org/10.5194/acp-11-1417-2011)
- 635 11-1417-2011, 2011.
- Mitchell, T. D.: Pattern Scaling. An Examination of the Accuracy of the Technique for Describing Future Climates, *Climatic Change*, 60, 217–242, <https://doi.org/10.1023/A:1026035305597>, 2003.
- Nychka, D., Hammerling, D., Krock, M., and Wiens, A.: Modeling and emulation of nonstationary Gaussian fields, *Spatial Statistics*, 28, 21–38, <https://doi.org/10.1016/j.spasta.2018.08.006>, 2018.
- 640 Olonscheck, D. and Notz, D.: Consistently estimating internal climate variability from climate model simulations, *Journal of Climate*, 30, 9555–9573, <https://doi.org/10.1175/JCLI-D-16-0428.1>, 2017.

- Osborn, T. J., Wallace, C. J., Harris, I. C., and Melvin, T. M.: Pattern scaling using ClimGen: monthly-resolution future climate scenarios including changes in the variability of precipitation, *Climatic Change*, 134, 353–369, <https://doi.org/10.1007/s10584-015-1509-9>, 2016.
- 645 Osborn, T. J., Wallace, C. J., Lowe, J. A., and Bernie, D.: Performance of pattern-scaled climate projections under high-end warming. Part I: Surface air temperature over land, *Journal of Climate*, 31, 5667–5680, <https://doi.org/10.1175/JCLI-D-17-0780.1>, 2018.
- Riahi, K., Rao, S., Krey, V., Cho, C., Chirkov, V., Fischer, G., Kindermann, G., Nakicenovic, N., and Rafaj, P.: RCP 8.5-A scenario of comparatively high greenhouse gas emissions, *Climatic Change*, 109, 33–57, <https://doi.org/10.1007/s10584-011-0149-y>, 2011.
- Rougier, J., Sexton, D. M. H., Murphy, J. M., and Stainforth, D.: Analyzing the climate sensitivity of the HadSM3 climate model using ensembles from different but related experiments, *Journal of Climate*, 22, 3540–3557, <https://doi.org/10.1175/2008JCLI2533.1>, 2009.
- 650 Sanderson, B. M., Knutti, R., and Caldwell, P.: Addressing interdependency in a multimodel ensemble by interpolation of model properties, *Journal of Climate*, 28, 5150–5170, <https://doi.org/10.1175/JCLI-D-14-00361.1>, 2015.
- Santer, B., Wigley, T., Schlesinger, M., and Mitchell, J.: Developing climate scenarios from equilibrium results. Max-Planck-Institut für Meteorologie report., Tech. Rep. 47, 1990.
- 655 Sato, M., Hansen, J. E., McCormick, M. P., and Pollack, J. B.: Stratospheric aerosol optical depths, 1850–1990, *Journal of Geophysical Research*, 98, 22 987–22 994, <https://doi.org/10.1029/93JD02553>, 1993.
- Scher, S.: Toward Data-Driven Weather and Climate Forecasting: Approximating a Simple General Circulation Model With Deep Learning, *Geophysical Research Letters*, 45, 12,616–12,622, <https://doi.org/10.1029/2018GL080704>, 2018.
- Scher, S. and Messori, G.: Weather and climate forecasting with neural networks: using GCMs with different complexity as study-ground, *Geoscientific Model Development*, 12, 2797–2809, <https://doi.org/10.5194/gmd-12-2797-2019>, 2019.
- 660 Seneviratne, S. I., Nicholls, N., Easterling, D., Goodess, C., Kanae, S., Kossin, J., Luo, Y., Marengo, J., McInnes, K., Rahimi, M., Reichstein, M., Sorteberg, A., Vera, C., and Zhang, X.: Changes in climate extremes and their impacts on the natural physical environment., in: *Managing the Risk of Extreme Events and Disasters to Advance Climate Change Adaptation. A Special Report of Working Groups I and II of the Intergovernmental Panel on Climate Change (IPCC)*, edited by Field, C. B., Barros, V., Stocker, T. F., Qin, D., Dokken, D. J., Ebi, K. L., Mastrandrea, M. D., Mach, K. J., Plattner, G.-K., Allen, S. K., Tignor, M., and Midgley, P. M., chap. 3, pp. 109–230, Cambridge University Press, Cambridge, UK, and New York, NY, USA, 2012.
- 665 Seneviratne, S. I., Donat, M. G., Pitman, A. J., Knutti, R., and Wilby, R. L.: Allowable CO₂ emissions based on regional and impact-related climate targets, *Nature*, 529, 477–483, <https://doi.org/10.1038/nature16542>, 2016.
- Serreze, M. C. and Barry, R. G.: Processes and impacts of Arctic amplification: A research synthesis, *Global and Planetary Change*, 77, 85–96, <https://doi.org/10.1016/j.gloplacha.2011.03.004>, 2011.
- 670 Sutton, R. T., Dong, B., and Gregory, J. M.: Land/sea warming ratio in response to climate change: IPCC AR4 model results and comparison with observations, *Geophysical Research Letters*, 34, L02 701, <https://doi.org/10.1029/2006GL028164>, 2007.
- Taylor, K. E., Stouffer, R. J., and Meehl, G. A.: An overview of CMIP5 and the experiment design, *Bulletin of the American Meteorological Society*, 93, 485–498, <https://doi.org/10.1175/BAMS-D-11-00094.1>, 2012.
- 675 Tebaldi, C. and Arblaster, J. M.: Pattern scaling: Its strengths and limitations, and an update on the latest model simulations, *Climatic Change*, 122, 459–471, <https://doi.org/10.1007/s10584-013-1032-9>, 2014.
- Tebaldi, C. and Knutti, R.: Evaluating the accuracy of climate change pattern emulation for low warming targets, *Environmental Research Letters*, 13, 55 006, <https://doi.org/10.1088/1748-9326/aabef2>, 2018.
- Trenberth, K. E.: The Definition of El Niño, *Bulletin of the American Meteorological Society*, 78, 2771–2778, [https://doi.org/10.1175/1520-0477\(1997\)078<2771:TDOENO>2.0.CO;2](https://doi.org/10.1175/1520-0477(1997)078<2771:TDOENO>2.0.CO;2), 1997.
- 680

- Wartenburger, R., Hirschi, M., Donat, M. G., Greve, P., Pitman, A. J., and Seneviratne, S. I.: Changes in regional climate extremes as a function of global mean temperature: An interactive plotting framework, *Geoscientific Model Development*, 10, 3609–3634, <https://doi.org/10.5194/gmd-10-3609-2017>, 2017.
- 685 Weigel, A. P.: Ensemble Forecasts, in: *Forecast Verification: A Practitioner’s Guide in Atmospheric Science*, edited by Jolliffe, I. T. and Stephenson, D. B., chap. 8, pp. 141–166, John Wiley & Sons, Chichester, UK, 2nd edn., <https://doi.org/10.1002/9781119960003.ch8>, 2012.
- Wild, M.: Enlightening global dimming and brightening, *Bulletin of the American Meteorological Society*, 93, 27–37, <https://doi.org/10.1175/BAMS-D-11-00074.1>, 2012.
- 690 Williamson, D., Goldstein, M., Allison, L., Blaker, A., Challenor, P., Jackson, L., and Yamazaki, K.: History matching for exploring and reducing climate model parameter space using observations and a large perturbed physics ensemble, *Climate Dynamics*, 41, 1703–1729, <https://doi.org/10.1007/s00382-013-1896-4>, 2013.