

ClimateBench v1.0: A benchmark for data-driven climate projections

D. Watson-Parris¹, Y. Rao², D. Olivié³, Ø. Seland³, P. Nowack⁴, G. Camps-Valls⁵, P. Stier¹, S. Bouabid⁶, M. Dewey⁷, E. Fons⁸, J. Gonzalez⁹, P. Harder^{1,10}, K. Jeggle⁸, J. Lenhardt⁹, P. Manshausen¹, M. Novitasari¹¹, L. Ricard¹², C. Roesch¹³

¹Atmospheric, Oceanic and Planetary Physics, Department of Physics, University of Oxford, Oxford, UK

²North Carolina Institute for Climate Studies, North Carolina State University, Asheville, NC 28801, USA

³Norwegian Meteorological Institute, Oslo, Norway

⁴Climatic Research Unit, School of Environmental Sciences, Norwich, UK

⁵Image Processing Laboratory, Universitat de València, València, Spain

⁶Department of Statistics, University of Oxford, Oxford, UK

⁷Department of Meteorology, Stockholm University, Stockholm, Sweden

⁸Institute of Atmospheric and Climate Science, ETH Zurich, 8092 Zurich, Switzerland

⁹Institute for Meteorology, Universität Leipzig, Leipzig, Germany

¹⁰Fraunhofer ITWM, Kaiserslautern, Germany

¹¹Department of Electronic and Electrical Engineering, University College London, London, UK

¹²Laboratory of Atmospheric Processes and their Impacts, School of Architecture, Civil & Environmental Engineering, Ecole Polytechnique Fédérale de Lausanne, Lausanne, Switzerland

¹³School of Geosciences, University of Edinburgh, Edinburgh, UK

Correspondence to: Duncan Watson-Parris (duncan.watson-parris@physics.ox.ac.uk)

Key Points:

- We introduce the first benchmark for emulation of key spatially resolved climate variables derived from a full complexity Earth System Model
- Three baseline emulators are presented which are able to predict regional temperature and precipitation with varying skill
- Evaluation metrics and areas for future research are presented to encourage further development of trustworthy data-driven climate emulators

30
31

Abstract

32 Many different emission pathways exist that are compatible with the Paris climate agreement, and many more are possible
33 that miss that target. While some of the most complex Earth System Models have simulated a small selection of Shared
34 Socioeconomic Pathways, it is impractical to use these expensive models to fully explore the space of possibilities. Such
35 explorations therefore mostly rely on one-dimensional impulse response models, or simple pattern scaling approaches to
36 approximate the physical climate response to a given scenario. Here we present ClimateBench - the first benchmarking
37 framework based on a suite of CMIP, AerChemMIP and DAMIP simulations performed by a full complexity Earth System
38 Model, and a set of baseline machine learning models that emulate its response to a variety of forcings. These emulators can
39 predict annual mean global distributions of temperature, diurnal temperature range and precipitation (including extreme
40 precipitation) given a wide range of emissions and concentrations of carbon dioxide, methane and aerosols, allowing them to
41 efficiently probe previously unexplored scenarios. We discuss the accuracy and interpretability of these emulators and
42 consider their robustness to physical constraints such as total energy conservation. Future opportunities incorporating such
43 physical constraints directly in the machine learning models and using the emulators for detection and attribution studies are
44 also discussed. This opens a wide range of opportunities to improve prediction, robustness and mathematical tractability. We
45 hope that by laying out the principles of climate model emulation with clear examples and metrics we encourage engagement
46 from statisticians and machine learning specialists keen to tackle this important and demanding challenge.

Plain Language Summary

48 Many different emission pathways exist that are compatible with the Paris climate agreement, and many more are possible
49 that miss that target. While some of the most complex Earth System Models have simulated a small selection of possible
50 futures, it is impractical to use these expensive models to fully explore the space of possibilities. Such explorations therefore
51 mostly rely on simple approximations of the global mean temperature response to a given scenario. Here we present
52 ClimateBench - the first benchmarking framework based on a suite of state-of-the-art simulations performed by a full
53 complexity Earth System Model, and a set of baseline machine learning models that emulate its response to a variety of
54 forcings. These emulators can predict annual mean global distributions of temperature, diurnal temperature range and
55 precipitation (including extreme precipitation) given a wide range of emissions and concentrations of carbon dioxide,
56 methane and aerosols, allowing them to efficiently probe previously unexplored scenarios. We also describe a set of
57 evaluation metrics which we hope will entice statisticians and machine learning experts to tackle this important and
58 demanding challenge.

59 1. **Introduction**

60 Many different emission pathways exist that are compatible with the Paris Agreement of limiting global mean temperatures
61 to “well below 2°C above pre-industrial levels and pursuing efforts to limit the temperature increase to 1.5 °C”, and many
62 more are possible that miss that target. Sampling possible emissions scenarios is therefore crucial for policy makers to weigh
63 the economic cost and societal impact of different mitigation and adaptation strategies. While many of the most complex
64 Earth System Models (ESMs) have simulated a small selection of ‘Shared Socioeconomic Pathways’ (SSPs; self-consistent
65 emissions scenarios based on assumptions about future socio-economic changes and imperatives) it is impractical to use
66 these expensive models to fully explore the space of possibilities (O’Neill et al., 2016). Therefore, such explorations mostly
67 rely on one-dimensional impulse response models, or simple pattern scaling approaches to approximate the physical climate
68 response to a given scenario (e.g., Millar et al., 2017).

69 Impulse response models (Smith et al., 2018; Meinshausen et al., 2011; Nicholls et al., 2020) are physically interpretable and
70 can capture the general non-linear behaviour of the system, but are inherently unable to model regional climate changes,
71 while pattern scaling approaches rely on a simple scaling of spatial distributions of temperature (e.g., Tebaldi et al., 2014) by
72 global mean temperature changes. This approach breaks down when considering precipitation, however, because of the
73 strong non-linearities in its response to temperature (e.g., Cabré et al., 2010). Statistical emulators of the regional climate
74 have been developed although these have been quite bespoke (Castruccio et al., 2014) or focus on the relatively simple
75 problem of emulating temperature (Holden and Edwards, 2010). These approaches also do not account for the influence of
76 aerosol, which can be important for both regional temperature and precipitation (e.g. Kasoar et al. 2018, Wilcox et al. 2020).
77 As has been noted recently (Watson-Parris, 2021), approaches including non-linear pattern scaling (Beusch et al., 2020) and
78 Gaussian process (GP) regression of long-term climate responses (Mansfield et al., 2020) suggest the possibility of using
79 modern machine learning (ML) tools to produce robust and general emulators of future scenarios. However, comparing and
80 contrasting these approaches is currently hindered by the lack of a consistent benchmark.

81 ClimateBench defines a set of criteria and metrics for objectively evaluating such climate model emulation; aims to
82 demonstrate the feasibility of such emulators; and provides a curated dataset that will allow, and hopefully encourage,
83 broader engagement with this challenge in the same way WeatherBench (Rasp et al., 2020) has achieved for weather
84 modelling. The target is to predict annual mean global distributions of temperature (T), diurnal temperature range (DTR),
85 precipitation (PR) and the 90th percentile of precipitation (PR90). These variables are chosen to represent a range of
86 important climate variables which respond differently to each forcing and include extreme changes (PR90) that might not be
87 expected to scale in the same way as the mean. For example, while T has been shown to scale roughly linearly with global
88 mean temperature changes (Castruccio et al., 2014), PR responds non-linearly, and DTR is more sensitive to aerosol
89 perturbations than global mean temperature changes (Hansen et al., 1995). Four of the main anthropogenic forcing agents are
90 provided as emulator inputs (predictors): carbon dioxide (CO₂), sulphur dioxide (SO₂; a precursor to sulphate aerosol),
91 black carbon (BC) and methane (CH₄). To enable spatially accurate emulators ClimateBench includes (annual mean): spatial

92 distributions of emissions for the short-lived aerosol species (SO₂ and BC), globally averaged emissions of CH₄, and global
93 cumulative emissions of CO₂.

94 The training data which is provided in order to support such predictions is generated from the simulations performed by the
95 second (and latest) version of the Norwegian Earth System Model (NorESM2; Seland et al., 2020) as part of the sixth
96 coupled model intercomparison project (CMIP6; Eyring et al., 2016). The provided inputs are constructed from the same
97 input data that is used to drive the original simulations. While we could have included simulations from multiple different
98 models, only one model submitted all of the DECK (Diagnostic, Evaluation, and Characterization of Klima), historical,
99 AerChemMIP (Collins et al., 2017) and ScenarioMIP (O'Neill et al., 2016) experiments required for our purposes, making it
100 impossible to provide a harmonised dataset. Further, there is no agreed way of robustly combining multiple models, and
101 while statistically combining multiple different models can lead to improved skill (Pincus et al., 2008) the resulting variance
102 is not reliable since the models are not truly independent (Knutti et al., 2013). Nevertheless, this single model dataset still
103 allows us to explore both scenario uncertainty and internal variability. Further, since even very simple models are able to
104 capture a variety of forcing responses (Smith et al. 2021), there is reason to believe that the response of the models to a given
105 forcing is more consistent than the range of responses (e.g., Richardson et al. 2019). We thus suggest that an emulator that
106 works best for NorESM2 will also have the tendency to perform better in emulating other CMIP models, mainly because the
107 data characteristics are by design similar (CMIP models represent the same physical system). In contrast, variations in the
108 structure of learning algorithms vary more significantly and follow entirely different ways of building a regression model.

109 As a demonstration of the variety of possible approaches to tackle this benchmark we also introduce three distinct baseline
110 emulators trained and evaluated against ClimateBench. These constitute the first data driven models for the projection of
111 multiple climatic variables and show promising skill in both the global-mean and spatial responses. We discuss the merits
112 and challenges in using each class of (regression) model and hope these provide a useful starting point for researchers
113 wishing to develop more advanced emulators.

114 The remainder of this paper describes the development of the dataset including the underlying ESM and all post-processing
115 (Section 2), the evaluation metrics used to rank ClimateBench submissions (Section 3), the baseline emulators (Section 4), a
116 discussion of such approaches and future opportunities for diverse approaches (Section 5) before providing a few concluding
117 remarks in Section 6.

118 2. Data set description and preparation

119 The data provided as part of ClimateBench is a heavily curated version of that publicly available in the CMIP6 data archive.
120 Here we describe the data extraction and processing steps, but the scripts used to perform this are also freely available (as
121 described in the data availability statement).

122 We use a selection of complementary simulations in order to provide as large a training dataset as possible while attempting
123 to avoid unnecessary redundancy. Table 1 details the full list of simulations included, the period they cover and a brief
124 description of their purpose in this context. Given that the primary purpose of ClimateBench is to train emulators over
125 different emission scenarios, ScenarioMIP simulations are a key component of the dataset. ScenarioMIP prescribes a limited
126 set of possible future emissions pathways exploring different socio-economic scenarios representing plausible narratives.
127 These scenarios are designed to span a range of mitigation scenarios (denoted by the first number in each scenario) and end-
128 of-century forcing possibilities (denoted by the last two numbers in each scenario). We include all available simulations,
129 including the AerChemMIP *ssp370-lowNTCF* variation of *ssp370* which includes lower emissions of near-term climate
130 forcers (NTCFs) such as aerosol (but not methane). We choose *ssp245* as our test dataset against which all ClimateBench
131 emulators are to be evaluated. This scenario represents a medium mitigation and medium forcing scenario, ensuring trained
132 emulators are able to interpolate a solution rather than extrapolate (as discussed further in Section 5). The CMIP6 *historical*
133 experiment is also included since it provides useful training data at low emissions values.

134

135 **Table 1: Details of post-processed simulations provided as part of the ClimateBench dataset. Experiments denoted (*) are ancillary**
136 **data that, while potentially useful, are not used in training the baseline emulators presented here.**

Protocol	Experiment	Period	Notes
ScenarioMIP (O'Neill et al., 2016)	ssp126	2015 - 2100	A high ambition scenario designed to produce significantly less than 2 degrees warming by 2100.
	ssp245	2015 - 2100	Designed to represent a medium forcing future scenario. This is the test scenario to be held back for evaluation.
	ssp370	2015 - 2100	A medium-high forcing scenario with high emissions of near-term climate forcers (NTCF) such as methane and aerosol.
	ssp370-lowNTCF	2015 - 2054	Variation of SSP370 with lower emissions of aerosol and their precursors
	ssp585	2015 - 2100	This scenario represents the high end of the range of future pathways in the IAM literature and leads to a very large forcing of 8.5 Wm^{-2} in 2100.
CMIP6 (Eyring et al., 2016)	historical	1850 – 2014	A simulation using historical emissions of all forcing agents designed to recreate the historically observed climate.
	abrupt-4xCO2*	500 years	Idealised simulation in which CO2 is abruptly quadrupled. Other forcing agents remain unchanged.
	1pctCO2*	150 years	Idealised simulation in which CO2 is gradually increased by 1% / year. Other forcing agents remain unchanged.
	piControl*	500 years	Baseline simulation in which all forcing agents remain unchanged.

DAMIP (Gillett et al., 2016)	hist-GHG	1850 – 2014	A historical simulation with varying concentrations for CO ₂ and other long-lived greenhouse-gases (only).
	hist-aer	1850 – 2014	A historical simulation only forced by changes in anthropogenic aerosol.
	ssp245-aer	2015 - 2100	A medium forcing scenario with only changes in anthropogenic aerosol, which provides an alternative test scenario for emulator evaluation.

137

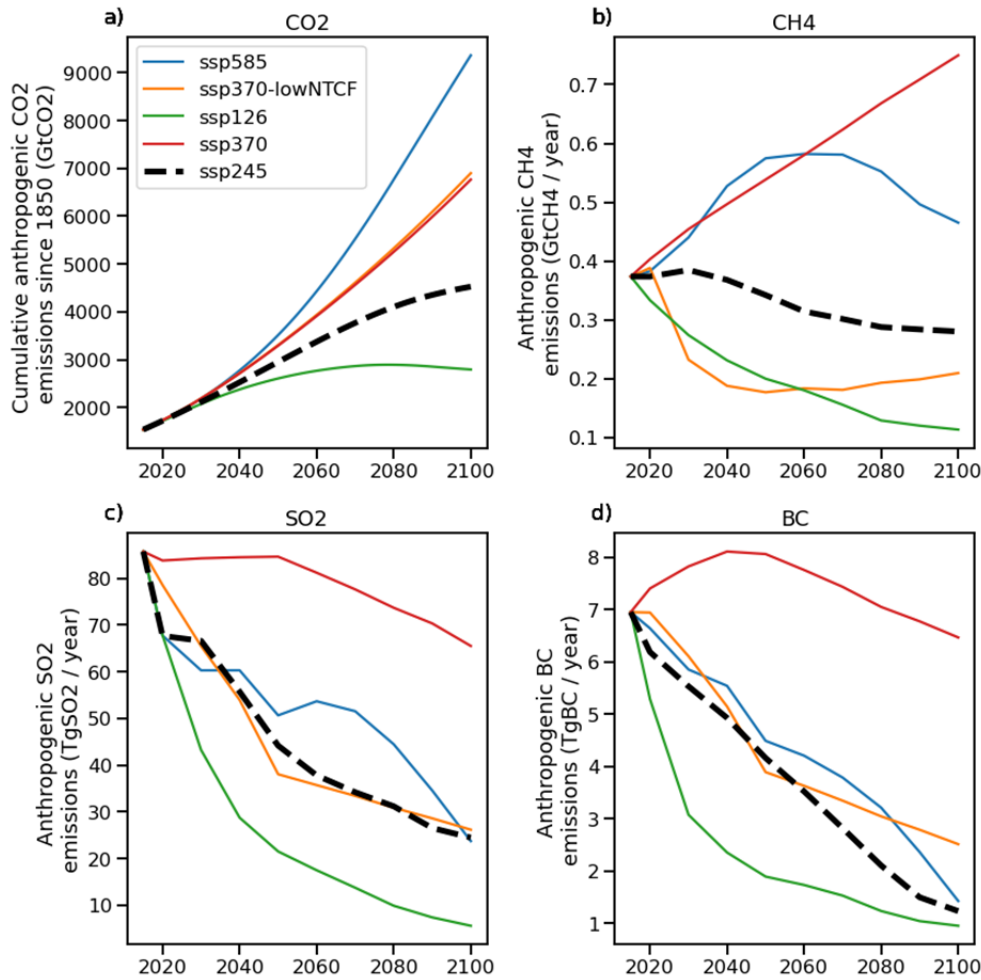
138 ClimateBench also includes a selection of more idealised simulations which are intended to provide training data at the
139 ‘corners’ of the four-dimensional input space, again helping reduce the chances of extrapolation in the resulting emulators
140 (as demonstrated in Figure A1). Two simulations that are commonly used to diagnose the equilibrium and transient climate
141 sensitivity are *abrupt-4xCO2* and *1pctCO2*, respectively. As the name suggests, the *abrupt-4xCO2* includes an abrupt
142 quadrupling of CO₂ over the pre-industrial concentrations while all other forcing agents remain unchanged. This level of
143 concentration represents the high end of future scenarios, broadly in line with *ssp585* but with no contribution from the other
144 forcers, simplifying its interpretation. The abrupt nature of the forcing also allows the timescale of the responses to be
145 determined which can be useful for emulators which account for this. The *1pctCO2* simulation gradually increases the
146 atmospheric concentration of CO₂ by 1% per year, again with other forcing agents unchanged. While potentially very useful,
147 they are not used in the training of the emulators presented in this work. Two other idealised simulations performed as part
148 of the Detection-Attribution Model Intercomparison Project (DAMIP; Gillett et al., 2016) represent the historical period
149 forced by only CO₂ and other long-lived greenhouse gases (*hist-GHG*), or only anthropogenic aerosol (*hist-aer*). These
150 provide opportunities to train emulators in regions of the input (emissions) space that are at the limits of plausible future
151 scenarios and were used in training the emulators described in Section 4.

152 Finally, the *piControl* simulation provides a baseline simulation with all forcings remaining unchanged from their pre-
153 industrial values. All target variables are calculated as a change against this climatology to simplify the training and
154 interpretation of the results. This long (500 year) simulation also enables a robust estimation of internal variability of the
155 climate system for those emulators which are able to represent it in future work, as discussed further in Section 5.1.

156 2.1 Input variables

157 The input data for these simulations is prescribed by the experimental protocol and provided by the input4MIPS project
158 (<https://esgf-node.llnl.gov/search/input4mips/>), which we collate and pre-process for ease of use. Specifically, we extract the
159 provided global mean emissions of CO₂ and CH₄ for each of the realistic (historical, ScenarioMIP and DAMIP) experiments
160 from the checksum files provided by the Community Emissions Data System (CEDS) dataset (Hoesly et al., 2018). We sum
161 over each sector and each month in order to derive annual total emissions and convert from Kg to Gt of CO₂. Some
162 historical and future periods are only provided in 5-yearly increments, so we linearly interpolate to yearly values for
163 consistency. The CO₂ emissions are summed cumulatively since, for realistic scenarios, a compensation between forcing

164 efficiency and ocean uptake means the temperature response to CO₂ is approximately linear in the cumulative emissions
 165 (Matthews and Caldeira 2008; Allen et al. 2009). Figure 1 shows the global mean emissions of each of the forcing agents
 166 under different future emissions scenarios, showing a wide range of possible pathways.
 167



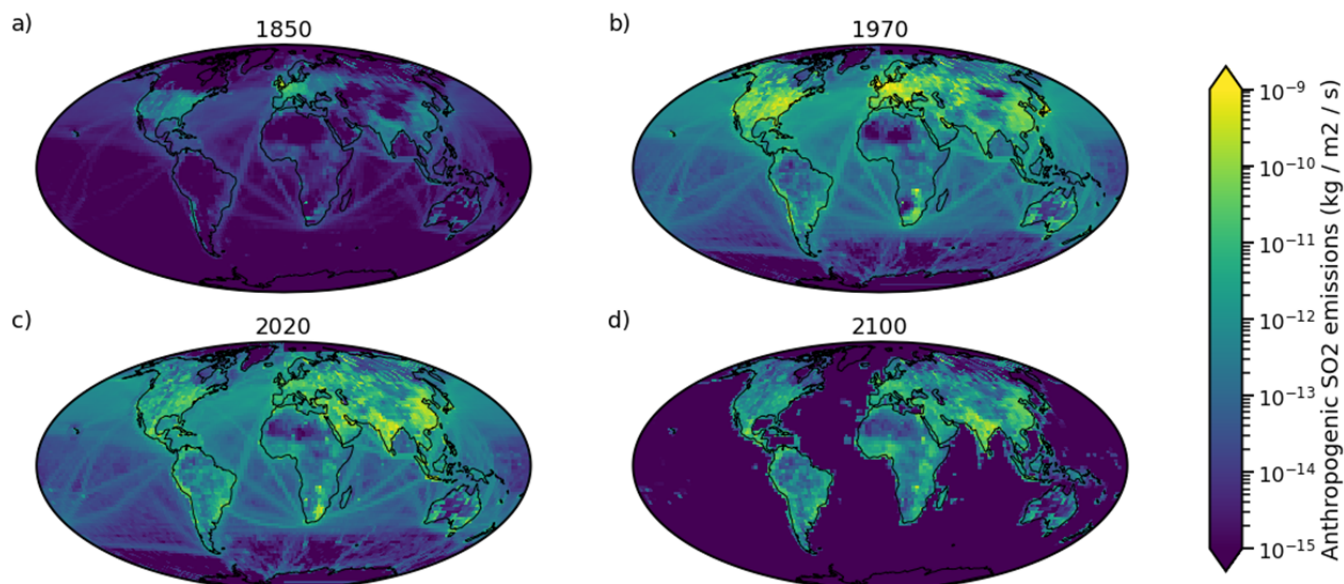
168
 169 **Figure 1: Time series of cumulative anthropogenic CO₂ emissions since 1850 (a); emissions of CH₄ (b); global mean emissions of**
 170 **SO₂ (c) and black carbon (BC; d) derived from NorESM2 ScenarioMIP simulations available within ClimateBench, including the**
 171 **SSP245 test scenario (shown in black).**

172
 173 The aerosol (precursor) emissions are derived from the latest version of the spatially resolved CEDS dataset and again
 174 summed over sectors and months to produce maps of annual total emissions, as shown in Figure 2 for SO₂ in different years.
 175 While the spatial distribution clearly evolves over the historical period and into the future scenarios, the emissions are fairly

176 localised around industrialised regions and dimensionality reduction can be used to reduce the size of these input features (as
177 discussed for the baseline emulators in Section 4). An area preserving interpolation is performed so that the emission data are
178 provided on the same spatial grid as the NorESM2 output fields to simplify its use in ML workflows. Again, as used for
179 NorESM2 the 5-yearly data is interpolated to a yearly frequency for consistency.

180

181



182

183 **Figure 2: Maps showing the evolution of the spatial distribution of anthropogenic SO₂ emissions in the pre-industrial era**
184 **represented by 1850 (a); the peak emissions era of 1970 (b); current emissions (c); and future emissions under SSP 245 (d).**

185 For the idealised CMIP simulations (*abrupt-4xCO₂* and *1pctCO₂*) no emissions files are used and so the cumulative
186 anthropogenic CO₂ emissions are calculated from the difference in the diagnosed CO₂ atmospheric mass concentrations in
187 these and the *piControl* experiment. Emissions of all other species are also provided but set to zero (as they represent no
188 change since the pre-industrial).

189 2.2 Target ESM

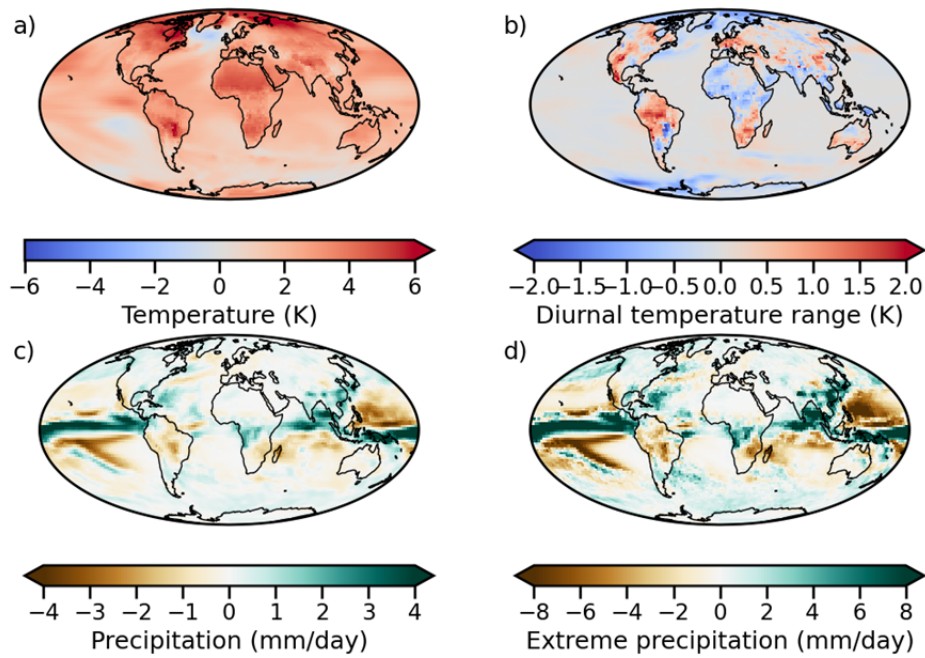
190 We use the output from simulations performed by the NorESM2 model in its low atmosphere-medium ocean resolution
191 (LM) configuration (Seland et al., 2020). This model consists of a fully coupled earth system with online atmosphere, land,
192 ocean, ice and biogeochemistry components. It shares many components with the Community Earth System Model Version
193 2 (Danabasoglu et al., 2020) but has a replaced aerosol and atmospheric chemistry scheme (including their interactions with
194 clouds) and a different ocean model. It has a relatively low equilibrium climate sensitivity (ECS; equilibrium global mean

195 temperature after a doubling of CO₂) of 2.5 K, particularly compared to the 5.3K of CESM2 (Gettelman et al., 2019), which
196 has been attributed to ocean heat uptake and convective mixing in the Southern Ocean (Gjermundsen et al., 2021). This is,
197 nevertheless, well within the assessed likely range of ECS (90% probability between 2–5°C; Forster et al., 2021) and makes
198 the emulation task harder than it might be for other CMIP6 models since the warming signal is weaker at the end of SSP245
199 (by which point CO₂ is the dominant forcing), as shown in Fig. A4. The combination of a weak ECS with a relatively strong
200 aerosol forcing (-1.36 W m^{-2} for 1850 to 2014), likely accounts for the somewhat anomalous cooling between 1950-1980 in
201 the historical simulations (Seland et al., 2020), although it has been noted that the combined anthropogenic response in
202 NorESM is realistic (Gillett et al. 2021).

203 **2.3 Output variables**

204 The output of these simulations are aggregated to annual mean values but kept at their native spatial resolution
205 (approximately 2°). The temperature (T) and precipitation (P) are exactly equivalent to the archived surface air temperature
206 (tas) and total precipitation (pr) output variables respectively. The DTR is calculated as the annual mean of the difference in
207 daily maximum and minimum surface air temperatures: $|tas_{max} - tas_{min}|_{ann}$. The PR90 is calculated as the 90th
208 percentile of the daily precipitation in each year. The annual mean baseline values (from the full *piControl* simulation) for
209 each variable are then subtracted from each experiment so that they represent a difference from pre-industrial. Temperature
210 changes under anthropogenic climate change are routinely reported in this way, and it also makes the downstream emulation
211 task somewhat easier as it removes an offset. The values are not scaled to have unit variance, but users of the dataset might
212 choose to do this with certain emulators. Many of the NorESM2 simulations include three ensemble members sampling
213 internal variability by choosing different initial model states from the start of the *piControl* simulation at intervals of 30
214 model years apart. These are included to allow (optional) emulation of internal variability.

215 Samples of these output fields from the target ssp245 dataset are shown in Figure 3. The relative increase in warming in the
216 northern polar regions (known as Arctic amplification) is clearly seen in Fig. 3a, as well as the north Atlantic warming hole
217 (Woollings et al., 2012; Drijfhout et al., 2012; Manabe and Stouffer, 1993), the emergence of which is also affected by
218 aerosol radiative forcing (Dagan et al., 2020). Figure 3b shows the strong land/sea contrast in DTR, since most of the change
219 is confined to land, and largely caused by changes in aerosol (particularly sulfate) forcing. Most of the precipitation response
220 shown in Figure 3c-d is due to the shift in the inter-tropical convergence zone (ITCZ) which results from a shift in the cross-
221 equatorial energy balance under increased warming (Schneider et al., 2014), but some features, particularly in South-East
222 Asia might be due to local aerosol effects (particularly due to BC; e.g., Bollasina et al. 2014, Wilcox et al. 2020, Mansfield
223 et al. 2020).



224

225 **Figure 3: Maps of target outputs from the SSP245 held-back test scenario at 2100 (as an anomaly to the pre-industrial control run)**
 226 **performed by NorESM2: (a) Annual mean surface temperature; (b) annual mean diurnal surface temperature range; (c) annual**
 227 **mean precipitation; and (d) 90th percentile of the daily precipitation.**

228 Also included in the dataset are the top-of-atmosphere Effective Radiative Forcings (ERFs) for this model for each forcing
 229 agent over the historical period. These are based on diagnostics of the fixed sea-surface temperature experiments of the
 230 Radiative Forcing Model Intercomparison Project (RFMIP; Pincus et al., 2016; Smith et al., 2021) and provide a more direct
 231 estimate of the radiative climate effect of each forcer over this period than simply emissions. It also allows an estimate of the
 232 efficacy of each forcer in this model (the temperature response per unit of forcing). This might be useful for normalising the
 233 inputs by their efficacy or developing more physically interpretable emulators that derive the climate response via the
 234 forcing, but these are not used in the present study.

235 3. Benchmark task

236 The task defined by ClimateBench is the prediction of the output variables described in Section 2.3 using only the inputs
 237 available from Section 2.1 under the chosen test scenario - ssp245. Emulators may choose to use as much or as little of the
 238 data presented in Table 1 in order to train their models as appropriate for a given approach. They may also choose to predict
 239 the contemporaneous response to emissions (as used in our RF and GP baseline emulators), account for a lagged response (as
 240 in our baseline NN emulator), or even predict the full time-series simultaneously.

241 3.1. Evaluation metrics

242 The evaluation criteria are a crucial aspect to any benchmark dataset and need to be concretely defined and accurately reflect
 243 the objectives of the machine learning task. Ideally, the criteria are also simple to implement such that they can be used as a
 244 target in any loss function that might be used to train emulators. The global mean changes in temperature and precipitation
 245 are key climatic variables but the spatial characteristics of the outputs in this task also need to be considered if the emulators
 246 are to be used for regional projections. As a primary metric we choose to combine the normalised, global mean root-mean
 247 square error ($NRMSE_s$) and the NRMSE in the global mean ($NRMSE_g$), calculated following:

$$NRMSE_s = \sqrt{\langle (|x_{i,j,t}|_t - |y_{i,j,t,n}|_n)^2 \rangle} / \langle |y_{i,j}|_t \rangle_{t,n} \quad (1)$$

$$NRMSE_g = \sqrt{\langle (\langle x_{i,j,t} \rangle - \langle |y_{i,j,t,n}|_n \rangle)^2 \rangle} / \langle |y_{i,j}|_t \rangle_{t,n} \quad (2)$$

$$NRMSE_t = NRMSE_s + \alpha * NRMSE_g, \quad (3)$$

248 where the global mean denoted $\langle x_{i,j} \rangle$ includes a weighting function that accounts for the decreasing grid-cell area towards
 249 the poles and is defined as: $\langle x_{i,j} \rangle = \frac{1}{N_{lat}N_{lon}} \sum_i^{N_{lat}} \sum_j^{N_{lon}} \cos(lat(i))x_{i,j}$, and α is a coefficient empirically chosen to be 5 so
 250 that each component provides roughly equal weight.

251 Combining these commonly used metrics in this way provides a single number summarising the mismatch between the
 252 predictions (x) and the target variables (y). By squaring the difference, the RMSE also weighs large discrepancies more
 253 heavily, penalising larger errors. We average the target variables over the three available ensemble members (n) and a
 254 relatively long period of the target scenario (2080 – 2100) in order to minimise the contribution of internal variability. We
 255 choose the final years of the century since the start of the *ssp245* is quite similar to some of the training scenarios. We
 256 normalise the RMSEs so that the metrics are broadly comparable across the target variables.

257 Estimates of this internal variability can be very valuable for climate projections however and since ClimateBench includes
 258 three ensemble members for each training dataset emulators are encouraged to include estimates of it if they are able. A
 259 natural extension of the RMSE for probabilistic estimates commonly used in weather forecasting is the Continuous Ranked
 260 Probability Score (CRPS):

$$CRPS = \int_{x=-\infty}^{x=\infty} (\langle F_{i,j,t}(x) \rangle - \langle F_{i,j,t}(y) \rangle)^2 dx, \quad (4)$$

261 where $F(x)$ and $F(y)$ are the cumulative distribution functions (CDFs) over the predicted and target ensembles respectively
 262 (Gneiting et al. 2005). This measures the area between the two CDFs so that smaller values are better and has the benefit of
 263 retaining a well-defined interpretation in the case of only a single target observation (whose CDF would be the Heaviside
 264 function). The CDFs can be approximated over finite ensembles using quadrature, or direct integration if the PDFs can be
 265 assumed to be Gaussian. It should be noted that the relatively low number of ensemble members available in ClimateBench

266 will likely underestimate full internal variability and a larger ensemble (e.g., 100 members in Rogers et al., 2021) should be
267 used for robust estimation. Indeed, the formulation above only includes variability in the global mean since such small
268 ensembles are unlikely to capture regional variability. Methods to calculate both metrics based on the *climpred* (Brady and
269 Spring, 2021) package are provided in the example notebooks included with the dataset. While this metric is not included in
270 the headline ranking of ClimateBench approaches, we include an example approach using GPs which is discussed in more
271 detail in Section 4.1.

272 **3.2. Baseline evaluation**

273 Before evaluating some baseline statistical emulators, it is useful to consider two cases with which we hope to place the skill
274 of the data-driven approaches in a broader context. The first is the internal variability of the NorESM2 target ensemble
275 which provides an upper bound on the predictability of the scenario in the presence of the natural variability of the Earth
276 system. This is estimated as the standard deviation across the three NorESM2 ensemble members in *ssp245*. In practice, the
277 emulators can (and do) outperform this baseline because they target the mean over all three ensemble members, reducing the
278 effect of internal variability. The second is a comparison against the inter-model spread encountered within CMIP6 for the
279 variables of interest which, despite (as discussed above) not providing a robust model uncertainty, represents a lower bound
280 on the accuracy we would like our emulators to achieve.

281 Additionally, we introduce a linear pattern scaling model which uses independent linear regressions of each of the output
282 variables at each model grid cell given the global mean temperature response to the emissions (e.g., Tebaldi et al., 2014).
283 This approach is somewhat simpler than the other data driven models since it assumes access to an accurate impulse
284 response (or box) model to determine the global mean temperature but provides a useful baseline. We train the regression
285 models using the same training output data as the other emulators (described in the next section) but the only input is the
286 global mean temperature. We assume this is available at prediction time as well so that this constitutes a ‘perfect’ pattern
287 scaling approach.

288 **4. Baseline emulators**

289 Three baseline emulators are developed to demonstrate various potential approaches to tackling the machine learning
290 problem this dataset provides. These are performed using the Earth System Emulator (ESEm; Watson-Parris et al., 2021) to
291 provide a simple interface for non-ML experts and permit sampling the emulators for potential use in detection and
292 attribution workflows (as discussed in the Section 5). All three emulators are trained using all the available training data: the
293 historical data; *ssp126*; *ssp370*; *ssp585*; and the historical data with aerosol (*hist-aer*) and greenhouse gas (*hist-GHG*)
294 forcings only, leading to 754 training / validation points (which are nevertheless not fully independent). More details on
295 emulator specific data pre-processing, training procedure and results are described in each of the following subsections.

296 The emulators all perform skilfully, as summarised in Table 2 and Figure 4. The emulators also show broadly similar biases,
 297 particularly for precipitation where they all slightly underestimate increases (decreases) in tropical (subtropical) rainfall in
 298 the western Pacific. They also tend to overpredict northern-hemisphere warming while underpredicting warming elsewhere.
 299 This might suggest that these particular changes are driven by different climate forcings or longer time-scale changes than
 300 modelled in this study. A direct comparison of the emulator predictions and NorESM is shown in Figure A2. Overall, the
 301 neural network emulator performs the best in predicting temperature and precipitation changes, while the Gaussian process
 302 emulator performs best at predicting changes in the diurnal temperature range.

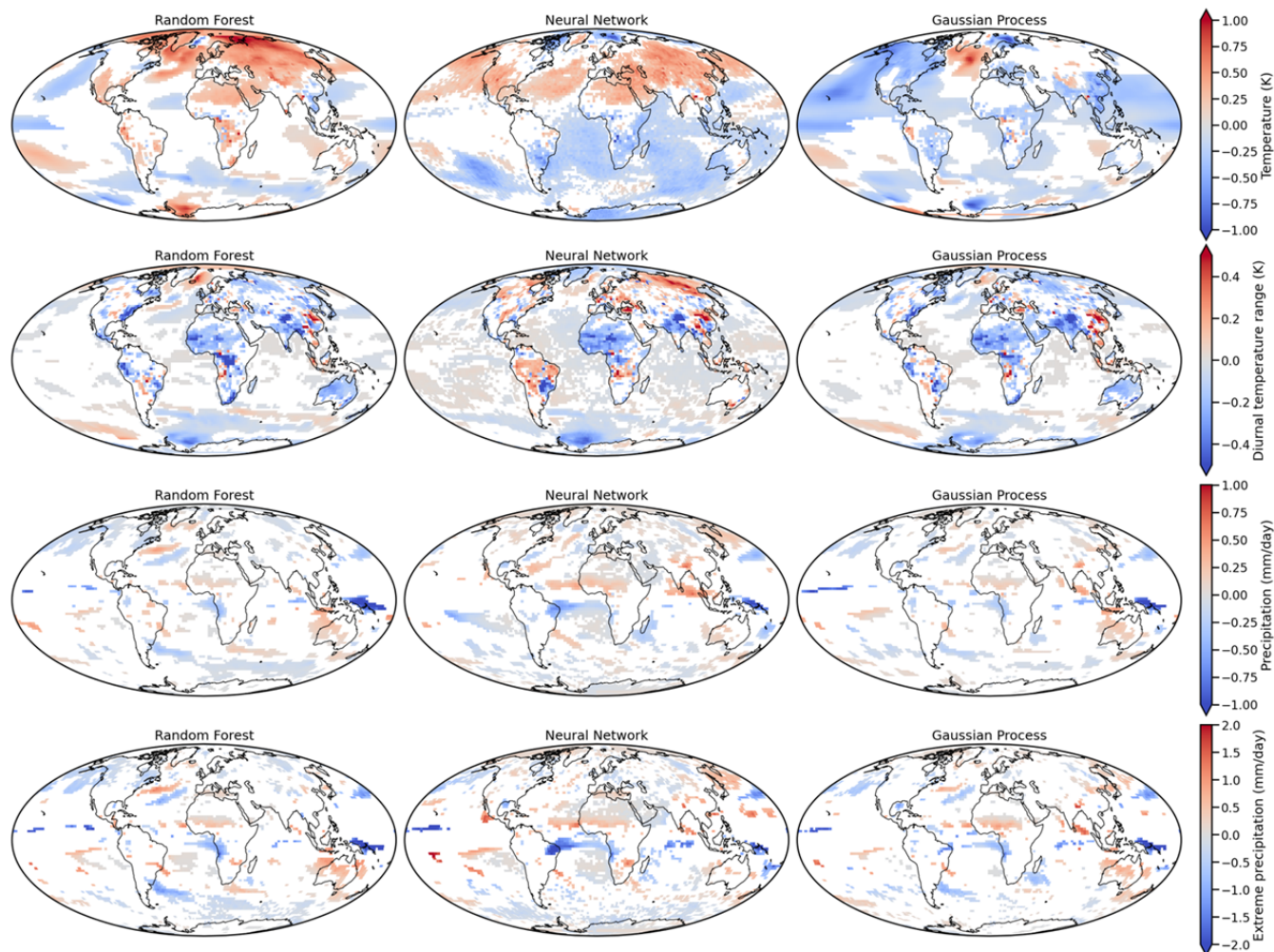
303

304 **Table 2: The spatial, global and total NRMSE of the different baseline emulators for the years 2080-2100 against the**
 305 **ClimateBench task of estimating key climate variables under future scenario SSP245. The best (lowest) emulator scores for each**
 306 **task are highlighted in bold. The normalised standard deviation in each variable over 22 different CMIP6 models and across the**
 307 **NorESM ensemble members are also included as indications of inter-model and internal variability, respectively.**

	NRMSE surface air temperature (1)			NRMSE diurnal temperature range (1)			NRMSE precipitation (1)			NRMSE 90 th percentile precipitation (1)		
	Spatial	Global	Total	Spatial	Global	Total	Spatial	Global	Total	Spatial	Global	Total
Gaussian Process	0.109	0.074	0.478	9.207	2.675	22.58 2	2.341	0.341	4.048	2.556	0.4 29	4.7 02
Neural Network	0.107	0.044	0.327	9.917	1.372	16.77 8	2.128	0.209	3.175	2.610	0.3 46	4.3 39
Random Forest	0.108	0.058	0.400	9.195	2.652	22.45 7	2.524	0.502	5.035	2.682	0.5 43	5.3 99
Pattern Scaling	0.080	0.048	0.320	8.083	2.327	19.71 9	2.006	0.331	3.662	2.400	0.4 12	4.4 61
Variability	0.052	0.072	0.414	2.513	1.492	9.973	1.350	0.268	2.691	1.757	0.4 57	4.0 43
CMIP6	-	0.206	-	-	0.815	-	-	0.417	-	-	-	-

308

309



311

312 **Figure 4: Maps of the mean difference in the ClimateBench target variables for each baseline emulator against the target NorESM**
 313 **values under the test ssp245 scenario averaged between 2080-2100. Differences insignificant at the $p < 5\%$ level are masked from**
 314 **the plots.**

315 4.1. Gaussian process regression

316 Gaussian processes (GPs) (Rasmussen and Williams, 2005) are probabilistic models which assume predictions can be
 317 modelled jointly as normally distributed. GPs have been widely used for nonlinear and nonparametric regression problems in
 318 the geosciences (Camps-Valls et al., 2016). A GP is fully determined by the expectation of individual predictions – referred
 319 to as the mean – and the covariance between pairs of predictions. Such covariance is typically user-specified as a bivariate
 320 function of the input data called the kernel function. The choice of the kernel function allows to restrict the functional class
 321 the GP belongs to, offering, for example, control over functional smoothness. GPs for regression solve a supervised problem

322 where the observed input-output sample pairs are used to: (1) infer the emulator parameters (typically only the noise variance
323 and the kernel parameters) by maximising the log-likelihood of the observations under the evidence; and then (2) allow to
324 obtain its posterior probability distribution that is used to make predictions over unseen inputs.

325 To prepare the input samples, the dimensionality of the SO₂ and BC emission maps are reduced with principal component
326 analysis, and we only use the 5 first principal components of each as inputs, corresponding to 96% and 98% of the explained
327 variance, respectively. All input covariates and target outputs are standardised using training data mean and standard
328 deviation.

329 The GP is set with a constant mean prior and separate kernels are devised for each species. Automatic relevance
330 determination (ARD) kernels are used for SO₂ and BC, allowing each principal component to be treated independently with
331 its own lengthscale parameter. The GP covariance function is obtained by summing all kernels together, thus accounting for
332 multiscale feature relations (see Camps-Valls et al., 2016 for several composite kernel constructions in remote sensing and
333 geoscience problems). To account for internal variability between ensemble members, we consider an additional white noise
334 term with constant variance over the output targets, which is also inferred from the training phase.

335 We use Matérn-1.5 kernels for each input. This guarantees the GP is a continuous, once differentiable function; details are
336 provided in Section A1. The mean value, kernels parameters and internal variability variance are jointly tuned against the
337 training data by marginal likelihood maximisation with the limited memory Broyden–Fletcher–Goldfarb–Shanno (L-BFGS)
338 optimisation algorithm. The emulators used have 18 parameters in total: 5 lengthscale and one variance parameter for each
339 aerosol kernel; one lengthscale and one variance for each of the GHG kernels; one mean and one likelihood variance.

340 As reported in Table 2, the total NRMSE of the mean predictions with the GPs are competitive with the neural network for
341 all the variables. This is remarkable given the limited number of parameters that are learnt. It suggests the GP prior is an
342 adequate choice for the purposes of emulation. Study of the inferred kernel variance (not shown) suggests that cumulative
343 CO₂ emissions generally influence all predictions, and unequivocally dominate the predictions for surface air temperature
344 and diurnal temperature range. CH₄ and BC emissions on the other hand appear to have negligible influence on the
345 predictions. Since the GP also provides posterior estimates of the variance (which will incorporate an estimate of internal
346 variability) we also calculate the CPRS for this emulator (see Table A5). While we are unable to compare these scores with
347 the other baseline methods the similarity to the global NRMSE indicates that the GP is also predicting the internal variability
348 accurately (otherwise it would be penalised in the CPRS relative to the NRMSE).

349 **4.2. Random Forests**

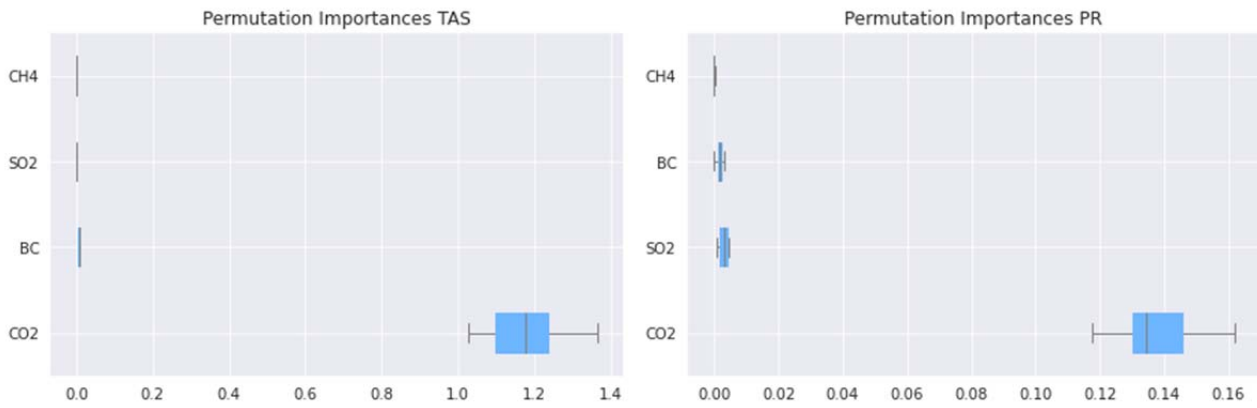
350 Random forests aggregate predictions of multiple decision trees (Ho, 1995; Breiman, 2001). These trees repeatedly split data
351 into subsets according to its features such that in-subset variance is low and between-subset variance is high. This makes
352 decision trees good at modelling non-linear functions, in particular interactions between different variables. However, they

353 are prone to overfitting (Ho, 1995). This problem is alleviated by ensemble methods which train a large number of different
354 trees. Weak learners are combined to give strong learners. Bagging, used in Random Forests, describes training different
355 trees on different subsets of the data or holding back some of the data dimensions for each individual tree. The Forest makes
356 a prediction by averaging over the predictions of all individual trees.

357 Two main arguments support an ensemble method approach to climate model emulation: These methods are skillful at
358 interpolation tasks, but by construction are unable to extrapolate (Breiman, 2001). However, for applications of climate
359 model emulation, interesting predictions will likely lie inside the hypercube delimited by historical data, low-emissions
360 (*ssp126*) and business-as-usual (*ssp585*) scenarios. A major advantage of ensemble methods over more complex ML
361 methods such as neural networks (and even ESMs) is their interpretability. This is important as ultimately predictions should
362 inform decision-making. Being able to provide explanations why a given input led to a prediction helps to understand the
363 consequences of decisions about emission pathways.

364 Analogously to the GP emulator, the dimensionality of aerosol emission maps is reduced with principal component analysis.
365 The first five principal components of SO₂ and BC together with the global emission maps of CO₂ and CH₄ form the input
366 features of the model. Separate random forest emulators are trained for the four target variables. The following
367 hyperparameters are tuned using random search of the training data without replacement: number of trees, tree depth,
368 number of samples required to split a node and to be at each leaf node. The hyperparameters used for each emulator are
369 indicated in Section A2.

370 As shown in Table 2, the spatial NRMSE scores of the random forest regressors are comparable to the performance of the
371 other emulators for all variables but the global NRMSE is significantly worse for temperature and precipitation (as can also
372 be seen in Figure 6). Discontinuities in the predicted global mean temperature change time series over this period (not
373 shown) perhaps indicate a deeper tree structure is required. To assess the impact of the four input features on the prediction,
374 we calculate the permutation feature importance. It is defined as the decrease in a model score when a single feature value is
375 randomly shuffled (Breiman, 2001). Figure 5 shows that CO₂ concentrations dominate the predictions. For temperature
376 predictions the other features are negligible. SO₂ and BC aerosol emissions have a small impact on the global mean
377 temperature and precipitation predictions. This is in line with the physical understanding that while anthropogenic aerosol
378 can influence precipitation rates (both radiatively and through aerosol-cloud interactions), aerosol contributions play a
379 negligible role at the end of the century in the *ssp245* test scenario. The regional influences may be more significant however
380 and this will be explored separately.



381

382

383

384

385

Figure 5: Permutation importances for the most important component of each variable in predicting global mean temperature (TAS) and precipitation (PR). Each emulator input variable is shuffled in turn to determine the relative contribution to prediction skill. Note that these average estimates do not account for potential regional contributions which may be particularly relevant for aerosol.

386

4.3. Neural Networks

387

388

389

390

391

392

393

394

395

396

Artificial Neural Networks (ANNs) are algorithms inspired by the biological neural networks of human brains that have shown outstanding success in areas like Computer Vision and Natural Language Processing. Two major ANN architectures are Convolutional Neural Networks (CNNs) (LeCun et al., 1990), that are able to model spatial dependencies, and Recurrent Neural Networks (RNNs), that are able to process time series and sequential data. ANNs have recently been employed to tackle a variety of problems in earth system science (Camp-Valls et al., 2021). CNNs are helpful for modeling climate data with a spatial structure, for instance, precipitation patterns or satellite imagery, and are frequently applied in climate science and weather forecasting (Trebing et al., 2020, Harder et al., 2020). Long short-term memory (LSTM) networks (Hochreiter et al., 1997), an advanced type of RNNs, have proven skillful for modeling climate time series, for example, for the prediction El Niño-Southern Oscillation (Broni-Bedaiko et al., 2019).

397

398

399

400

401

402

403

404

For time series of spatial variables, as in the ClimateBench dataset, we can use the two types of networks in sequence to model both spatial and temporal dependencies. The chosen architecture consists of a CNN followed by an LSTM built with the Keras library. The CNN includes one convolutional layer with 20 filters, a filter size of 3, and a ReLU activation function. The 3x3 pixel filters scan the input images to detect spatial patterns and feed these patterns to the next layer. These next layers are average pooling layers that reduce the spatial dimensionality ahead of the LSTM layer. The LSTM uses 25 units (i.e., the output dimension of each LSTM cell) and a ReLU activation function. The LSTM is followed by a dense layer and reshaping layer to (96, 144), i.e., the (latitude, longitude) dimension of the output variables.

405

406

407

408

409

The training data time-series is segmented into 10-year chunks, using a moving-time window in one-year increments, leading to 754 training samples of shape (10, 96, 144, 4) corresponding to the number of years, latitude, longitude and then number of variables. We trained four different emulators for the four different output variables. Each emulator is trained for 30 epochs, using a batch size of 16. For this baseline approach, we chose not to do any hyperparameter optimization, and all the parameters were chosen manually.

410
411 RMSE scores obtained with the CNN-LSTM architecture are somewhat better than those achieved with the other methods,
412 particularly in the global-mean. This might be because the LSTM is able to better capture the temporal autocorrelation than
413 the other emulators which treat the prediction instantaneously. The CNN-LSTM architecture also captures spatial changes in
414 temperature well (e.g., the Arctic amplification), even though warming at the poles is somewhat underestimated. In general,
415 warming in the Northern hemisphere is overestimated while it is underestimated in the Southern Hemisphere. Given the
416 overestimated temperature response in the *ssp245-aer* simulations shown in Figure A3, this may be due to an overestimation
417 of the effect of aerosol on the temperature by this emulator. The diurnal temperature range is well predicted, with a lower
418 performance over land. The CNN-LSTM also captures spatio-temporal changes in precipitation (e.g., the ICTZ shift) quite
419 well.
420

421 5. Discussion

422 5.1. Climate-specific challenges

423 The emulation of future climate states presents particular challenges for machine learning and other statistical approaches.
424 Chiefly among those is the limited amount of training data that is typically available; current ML approaches are not
425 prepared to learn such complex scenarios in small data regimes under a covariate shift. As pointed out, the complex ESMs
426 that are trusted to model the future climate are extremely computationally expensive to run and the observational record
427 cannot inform us about unseen future scenarios. By harnessing a large selection of simulations performed as part of CMIP6,
428 ClimateBench attempts to alleviate this difficulty, but nevertheless only around 500 training points (years) represent realistic
429 climate states, many of which are not independent (as shown in Fig. A1). This presents a challenge for deep learning
430 approaches which typically require tens of thousands of training samples to avoid over-fitting. The inclusion of longer
431 idealised simulations does provide opportunities for pre-training however, particularly the 500-year long *piControl*
432 simulations which could be used with contrastive learning to reduce the training samples required for neural network
433 architectures.

434 The *piControl* simulation could also be used to inform emulators more explicitly about the internal variability of climate (as
435 produced by NorESM2). The signal, particularly for the precipitation target variables, can be small compared to this
436 variability and this proves challenging for some emulators to reproduce. An explicit model of the internal variability
437 (Castruccio et al., 2019) could help to alleviate this.

438 Another challenge in applying statistical learning approaches to this dataset is the relatively high dimensional inputs and
439 outputs (96 x 144). Most approaches to emulating the regional temperature response to a CO₂ forcing have been carried out
440 at, at most, dozens of locations, but accounting for the spatial correlations is something which CNNs can excel at and have
441 recently been shown to produce accurate emulations of temperature across similar dimensionality (Beusch et al., 2020). Such
442 approaches typically assume a regular spacing, however, and neglect the reducing area of each grid-cell towards the poles.

443 While more traditional approaches of dimensionality reduction can also be used, such as (weighted) empirical orthogonal
444 functions (EOFs), these may not be appropriate for the non-linear precipitation fields which might require kernel-based
445 approximations (e.g., Bueso et al., 2019).

446 For practical purposes, an estimate of the uncertainty in any prediction would be extremely valuable. This uncertainty should
447 encompass that due to the internal variability and the emulator approximation (and ideally that of the underlying physical
448 model). In the ML community, these are known as the epistemic and the model uncertainties, and are being studied
449 intensively (Kendall et al., 2017). Quantifying these two uncertainties would allow increased trust (a concept explored in the
450 next section) in the prediction as well as quantitative comparison to other predictions. We encourage the estimation of
451 uncertainty wherever possible, using the provided CRPS metric to evaluate such probabilistic projections. The ability to
452 sample from such distributions would also permit the generation of so-called ‘superensembles’ which can provide very large
453 ensembles of multiple models under given scenarios (Beusch et al., 2020).

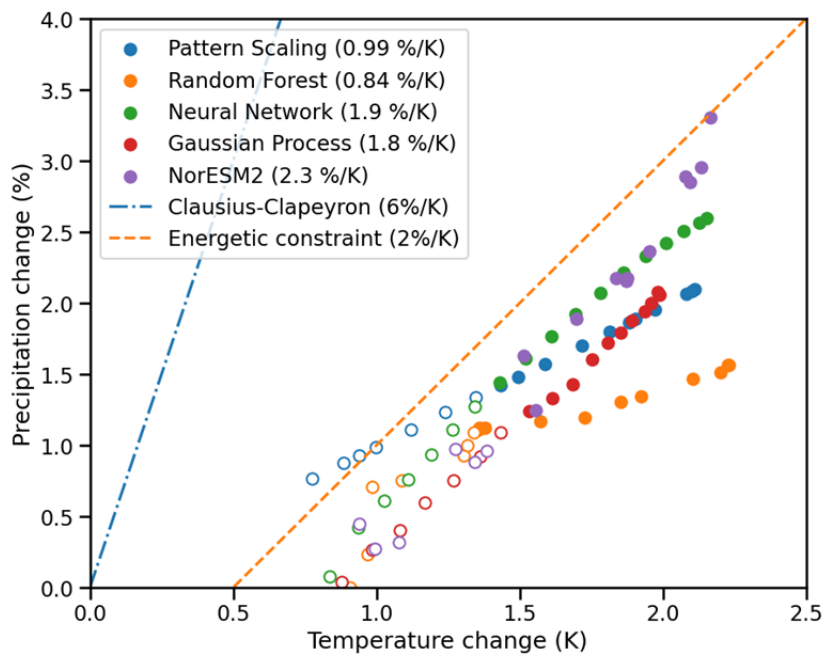
454 As previously discussed, and shown in Figure A4, there is large inter-model variability in the projected climate variables in
455 CMIP6, even across a single scenario. Future work should explore the ability of a given emulator to robustly recreate each of
456 these model responses, and could allow a deeper understanding of their discrepancies.

457 **5.2. Emulator trustworthiness**

458 For climate model emulators to be useful for policy decisions they must be trusted by their users. The trustworthiness of any
459 model is a subjective concept that broadly represents one's belief that the model faithfully represents some underlying ‘truth’.
460 Model verification attempts to objectively assert this view (indeed the word derives from the Latin, *verus*, meaning true) but
461 is formally impossible for an open system like the Earth (see e.g., Oreskes et al., 1994). While weather models can be
462 regularly validated against observations, in the climate sciences we often instead resort to necessarily incomplete model
463 evaluation and rely on underlying physical principles to provide reassurances of broader validity. The ClimateBench
464 emulators side-step this issue by aiming only to accurately reproduce an existing physical model which is assumed to already
465 be well evaluated, and therefore attain trustworthiness through proxy. It would nevertheless be reassuring if the emulators
466 could be demonstrated to respect some of the same physical constraints.

467 In this spirit, Figure 6 shows the relative change in global mean precipitation as a function of global mean temperature
468 change (the hydrological sensitivity) of the baseline emulators and NorESM2. While locally precipitation can change in
469 accordance with the Clausius-Clapeyron relationship (6-7% / K), energy conservation requires that the global changes in
470 precipitation are balanced by radiative cooling and limited to 2-3% / K (Allen and Ingram, 2002; Pendergrass and Hartmann,
471 2013; Jeevanjee and Romps, 2018; Dagan et al., 2019). While the RF emulator underestimates the hydrological sensitivity of
472 NorESM, it is clear that the emulators learn the physical relationship from the underlying model. Since the emulators were
473 trained on the precipitation and temperature this is to be expected to some degree, but this demonstrates the principle that
474 emulators trained correctly can retain the physical laws of the underlying models over the range of their training data. Future

475 efforts to introduce these invariances directly have the potential to significantly ease the training and improve the inference
476 of climate model emulators (Beucler et al., 2021), ultimately improving their trustworthiness.



477

478 **Figure 6: The relative change in global mean precipitation as a function of global mean temperature change in the**
479 **baseline emulators and NorESM2 averaged in 5-year increments to reduce internal-variability. Hollow and solid**
480 **points indicate years before and after 2050 respectively. The change predicted by the Clausius-Clapeyron**
481 **relationship and energy conservation considerations are shown as dashed lines.**

482 There has been much attention recently given to ‘interpretable’ and ‘explainable’ machine learning models, the former of
483 which are said to behave in a-priori understandable ways (Barnes et al., 2020), while the latter provide mechanisms to
484 determine post-hoc understanding (McGovern et al., 2019). While not as robust as physical laws, these techniques provide
485 useful indications that such models are getting the right answer for the right reasons. Indeed, the physical ESMs currently
486 considered the ‘gold standard’ of climate modelling are often only interpretable or explainable by expert practitioners and it
487 is hoped that (interpretable) ClimateBench emulators will be useful in analysing and understanding the response of the
488 underlying physical models themselves.

489 5.3. Research opportunities

490 While the challenges outlined above are mostly surmountable with modern architectures and carefully chosen workflows,
491 there are also several broader opportunities ClimateBench presents to develop the state-of-the-art in climate model
492 emulation.

493 As already mentioned, one area of particular interest is the use of hybrid modelling whereby statistical or ML based
494 emulators embed physical equations, constraints or symmetries in order to improve accuracy, robustness and generalisability
495 (Camps-Valls et al., 2021; Reichstein et al., 2019; Karpatne et al., 2017). One obvious way in which to apply such
496 approaches to ClimateBench is to marry the simple impulse response models discussed in Section 1 with more complex
497 methods to predict the spatial response. Such an approach has recently been demonstrated for temperature (Beusch et al.,
498 2021) but could conceivably be extended to modelling each of the fields targeted in ClimateBench. A more unified, and
499 ambitious, approach would be to model the ordinary differential equations of the response to a forcing directly in the
500 statistical emulator using either numerical GPs (Raissi et al., 2018) or Fourier neural operators (Li et al., 2020).

501 Another important open question when using data-driven approaches to emulate the climate is how to ensure predictions are
502 performed at locations within the distribution of the training data. In other words, how to ensure the emulator is being used
503 to interpolate existing model simulations rather than extrapolating to completely unseen regions of input space. This can be
504 easy to test for in low dimensions, but it becomes increasingly difficult in higher dimensions and while the training and test
505 data in ClimateBench have been chosen to minimise the risk of extrapolation broader use could be hindered by the risk of
506 inadvertently asking for an out-of-distribution prediction. While the predictive variance of GPs provide such indications (out
507 of the sample range the GP mean returns to the prior and the covariance is maximised), it is not so easy for other techniques
508 and the use of modern techniques to detect such occurrences (e.g., Lee et al., 2018; Rabanser et al., 2018) could be of great
509 value to minimise this risk.

510 **5.4. Application to detection and attribution**

511 The use of an efficient and accurate way of estimating the climate impacts of different emission scenarios is not limited to
512 exploring future pathways. We may also ask: ‘What observed climate states and events can be attributed to anthropogenic
513 emissions?’. A whole field, which started with the seminal work of Hasselmann (1993) has developed rapidly in the last
514 decade (Stott et al., 2016; Barnett et al., 2005; Stott et al., 2010; Shindell et al., 2009; Otto et al., 2016) attempting to answer
515 this question. A common approach is to use climate model (or ESM) simulations to determine optimal ‘fingerprints’ with
516 which to test observations as well as the power of such a fingerprint under internal variability. These typically have to make
517 fairly strong assumptions about the form of the climate response however (often relying on multiple linear regression) and
518 can incorporate observations of only a few dimensions.

519 One possible application of the efficient emulators trained using ClimateBench could then be to allow the inference of higher
520 dimensional attribution problems, incorporating more information (such as the DTR and PR) and potentially providing more
521 confident assessments. It would be straightforward to implement such an approach using the ESEm package which provides
522 a convenient interface for such inferences using e.g., approximate Bayesian computation, variational inference or Markov
523 Chain Monte-Carlo sampling. Future work will investigate these possibilities.

524 As a simple demonstration of potential of such an approach we have included a prediction by the emulators compared to the
525 original NorESM2 simulations of the *ssp245-aer* DAMIP experiment in which only the aerosol species are emitted, shown in
526 Figure A3. This is a more challenging scenario than the *ssp245* test case due to the much smaller total forcing and the
527 emulators do not perform as well (see NRMSE in Table A4). It is interesting to note that the emulators particularly struggle
528 with temperature changes in the North Atlantic where slow ocean circulation changes (e.g., Dagan et al. 2020) may not be
529 fully captured. They nevertheless capture the main features of the response and show promise for future work disentangling
530 the forcings and feedbacks in NorESM2, other ESMs and ultimately observations.

531 6. Conclusions

532 The application of machine learning to the prediction of future climate states has, perhaps justifiably due to the challenges
533 laid out above, been cautious to date. Particular applications however, with carefully chosen training data and objectives, can
534 provide fruitful avenues for research and open exciting opportunities for improvement over the current state-of-the-art. This
535 paper introduces the ClimateBench dataset in order to galvanise existing research in this area, provide a standard objective
536 with which to compare approaches and also introduce new researchers to the challenge of climate emulation. It provides a
537 diverse set of training data with clear objectives and challenging target variables, some of which have been extensively
538 studied (surface air temperature) and some which have been somewhat neglected (diurnal temperature range and
539 precipitation).

540 We also introduce three quite distinct approaches for undertaking this challenge: a random forest; a Gaussian process; and a
541 neural network model. These different models are based on different principles, have distinct assumptions and rely on quite
542 different learning paradigms. Each has their strengths and weaknesses but all perform well in the evaluation metrics and
543 generally reproduce the NorESM2 temperature and precipitation response well in a realistic (but unseen) future scenario,
544 especially compared to CMIP6 inter-model diversity. The neural network model performs best overall and shows good skill
545 both in the global mean and spatially. All the models perform less well in the aerosol only test, suggesting that they have not
546 fully learnt the distinct response due to each forcer and future emulators should aim to rectify this.

547 Current impact assessments are often based on simple emulators, which are then scaled to match modelled patterns, but
548 which are unable to predict non-linear responses in e.g., precipitation. A robust, trustworthy emulator which is able to
549 provide such predictions could be immensely valuable in quantifying and understanding the changes and associated risks of
550 different socio-economic pathways. Given the importance of faithfully and accurately reproducing the response of ESMs, we
551 hope the challenge will also spur innovation in nascent physically informed ML techniques.

552 In order to meet these objectives, we have provided open, easy to access datasets and training notebooks which reproduce
553 the results shown in this manuscript and demonstrate the use of the different baseline emulators. All software is open-source
554 and readily available using commonly used package managers. We hope this dataset will provide a focus for climate and ML

555 researchers to advance the field of climate model emulation and provide policy makers with the tools they require to make
556 well informed decisions.

557

558 **5. Data and code availability**

559 The baseline code is available on GitHub (<https://github.com/duncanwp/ClimateBench>) and a DOI for the specific version,
560 including that used to generate the plots in this paper, will be made available on acceptance.

561 The benchmark data is available here: <https://doi.org/10.5281/zenodo.5196512>. The raw CMIP6 data used here are available
562 through the Earth System Grid Federation and can be accessed through different international nodes e.g.: [https://esgf-
563 index1.ceda.ac.uk/search/cmip6-ceda/](https://esgf-index1.ceda.ac.uk/search/cmip6-ceda/).

564 **6. Author Contributions**

565 DWP conceptualised ClimateBench and performed the data-processing and initial analysis. YD and PN contributed to the
566 definition and setup of the framework. DO and ØS performed the original NorESM2 simulations used for training. SB, MD,
567 EF, PH, KJ, JL, PM, MN, LR, CR and JV developed the baseline emulators. DWP prepared the manuscript with
568 contributions from all co-authors.

569 **7. Acknowledgements**

570 DWP and PS acknowledge funding from NERC projects NE/P013406/1 (A-CURE) and NE/S005390/1 (ACRUISE). DWP,
571 GCV, PS, SB, MD, EF, PH, KJ, JL, PM, MN, LR, CR and JV acknowledge funding from the European Union's Horizon
572 2020 research and innovation programme iMIRACLI under Marie Skłodowska-Curie grant agreement No 860100. PS
573 additionally acknowledges support from the ERC project RECAP and the FORCeS project under the European Union's
574 Horizon 2020 research programme with grant agreements 724602 and 821205. GCV was partly supported by the European
575 Research Council (ERC) Synergy Grant "Understanding and Modelling the Earth System with Machine Learning
576 (USMILE)" under the Horizon 2020 research and innovation programme (Grant agreement No. 855187). YR was supported
577 by NOAA through the Cooperative Institute for Satellite Earth System Studies under Cooperative Agreement
578 NA19NES4320002. This research was supported, in part, by the National Science Foundation under Grant No. NSF PHY-
579 1748958.

580 The authors also gratefully acknowledge the World Climate Research Programme, which, through its Working Group on
581 Coupled Modelling, coordinated and promoted CMIP6. We thank the climate modelling groups for producing and making
582 available their model output, the Earth System Grid Federation (ESGF) for archiving the data and providing access, and the
583 multiple funding agencies who support CMIP6 and ESGF. In particular, DO and ØS acknowledge support from the Research
584 Council of Norway funded project INES (270061). High-performance computing and storage resources for NorESM2 were
585 provided by the Norwegian infrastructure for computational science (through projects NN2345K, NN9560K, NS2345K,

586 NS9560K, and NS9034K). We would also like to thank the many other participants of the 3rd NOAA AI Workshop
587 hackathon who provided valuable feedback.

588

589

590 8. References

591

592 Allen, M. R. and Ingram, W. J.: Constraints on future changes in climate and the hydrologic cycle, *Nature*, 419(6903), 224,
593 doi:10.1038/nature01092, 2002.

594 Allen, M. R., Frame, D. J., Huntingford, C., Jones, C. D., Lowe, J. A., Meinshausen, M. and Meinshausen, N.: Warming caused by
595 cumulative carbon emissions towards the trillionth tonne, *Nature*, 458(7242), 1163–1166, doi:10.1038/nature08019, 2009.

596 Barnes, E. A., Toms, B., Hurrell, J. W., Ebert-Uphoff, I., Anderson, C., and Anderson, D.: Indicator Patterns of Forced Change
597 Learned by an Artificial Neural Network, *J Adv Model Earth Sy*, 12, <https://doi.org/10.1029/2020ms002195>, 2020.

598 Barnett, T., Zwiers, F., Hengerl, G., Allen, M., Crowley, T., Gillett, N., Hasselmann, K., Jones, P., Santer, B., Schnur, R., Scott, P.,
599 Taylor, K., and Tett, S.: Detecting and Attributing External Influences on the Climate System: A Review of Recent Advances, *J*
600 *Climate*, 18, 1291–1314, <https://doi.org/10.1175/jcli3329.1>, 2005.

601 Beucler, T., Pritchard, M., Rasp, S., Ott, J., Baldi, P. and Gentine, P.: Enforcing Analytic Constraints in Neural Networks
602 Emulating Physical Systems, *Phys Rev Lett*, 126(9), 098302, doi:10.1103/physrevlett.126.098302, 2021.

603 Beusch, L., Gudmundsson, L., and Seneviratne, S. I.: Emulating Earth system model temperatures with MESMER: from global
604 mean temperature trajectories to grid-point-level realizations on land, *Earth Syst Dynam*, 11, 139–159, [https://doi.org/10.5194/esd-](https://doi.org/10.5194/esd-11-139-2020)
605 [11-139-2020](https://doi.org/10.5194/esd-11-139-2020), 2020.

606 Beusch, L., Nicholls, Z., Gudmundsson, L., Hauser, M., Meinshausen, M., and Seneviratne, S. I.: From emission scenarios to
607 spatially resolved projections with a chain of computationally efficient emulators: MAGICC (v7.5.1) – MESMER (v0.8.1)
608 coupling, *Geoscientific Model Dev Discuss*, 2021, 1–26, <https://doi.org/10.5194/gmd-2021-252>, 2021.

609 Bollasina, M. A., Ming, Y., Ramaswamy, V., Schwarzkopf, M. D. and Naik, V.: Contribution of local and remote anthropogenic
610 aerosols to the twentieth century weakening of the South Asian Monsoon: AEROSOLS AND SOUTH ASIAN MONSOON,
611 *Geophys Res Lett*, 41(2), 680–687, doi:10.1002/2013gl058183, 2014.

612 Brady, R. and Spring, A.: climpred: Verification of weather and climate forecasts, *J Open Source Softw*, 6, 2781,
613 <https://doi.org/10.21105/joss.02781>, 2021.

614 Breiman, L.: Random Forests, *Mach Learn*, 45(1), 5–32, doi:10.1023/a:1010933404324, 2001.

615 Broni-Bediako, Clifford & Katsriku, Ferdinand & Unemi, Tatsuo & Atsumi, Masayasu & Abdulai, Jamal-Deen & Shinomiya,
616 Norihiko & Owusu, Ebenezer Owusu. (2019). El Niño-Southern Oscillation forecasting using complex networks analysis of LSTM
617 neural networks. *Artificial Life and Robotics*. 24. 10.1007/s10015-019-00540-2.

618 Bueso, D., Piles, M. and Camps-Valls, G.: Nonlinear PCA for Spatio-Temporal Analysis of Earth Observation Data, *Ieee T Geosci*
619 *Remote*, 58(8), 5752–5763, doi:10.1109/tgrs.2020.2969813, 2020.

620 Cabré, M. F., Solman, S. A., and Nuñez, M. N.: Creating regional climate change scenarios over southern South America for the
621 2020's and 2050's using the pattern scaling technique: validity and limitations, *Climatic Change*, 98, 449–469,
622 <https://doi.org/10.1007/s10584-009-9737-5>, 2010.

623 Camps-Valls, G., Verrelst, J., Munoz-Mari, J., Laparra, V., Mateo-Jimenez, F., Gomez-Dans, J. and Gomez-Dan, J.: A Survey on
624 Gaussian Processes for Earth-Observation Data Analysis: A Comprehensive Investigation, *Ieee Geoscience Remote Sens Mag*,
625 4(2), 58–78, doi:10.1109/mgrs.2015.2510084, 2016.

626 Camp-Valls G., Tula D., Zhu X. X. and Reichstein M.: Deep Learning for the Earth Sciences: A Comprehensive Approach to
627 Remote Sensing, Climate Science and Geosciences, <https://onlinelibrary.wiley.com/doi/book/10.1002/9781119646181>, 2021

628 Castruccio, S., McInerney, D. J., Stein, M. L., Crouch, F. L., Jacob, R. L., and Moyer, E. J.: Statistical Emulation of Climate
629 Model Projections Based on Precomputed GCM Runs*, *J Climate*, 27, 1829–1844, <https://doi.org/10.1175/jcli-d-13-00099.1>, 2014.

630 Castruccio, S., Hu, Z., Sanderson, B., Karspeck, A., and Hammerling, D.: Reproducing Internal Variability with Few Ensemble
631 Runs Reproducing Internal Variability with Few Ensemble Runs, *J Climate*, 32, 8511–8522, [https://doi.org/10.1175/jcli-d-19-](https://doi.org/10.1175/jcli-d-19-0280.1)
632 0280.1, 2019.

633 Collins, W. J., Lamarque, J.-F., Schulz, M., Boucher, O., Eyring, V., Hegglin, M. I., Maycock, A., Myhre, G., Prather, M., Shindell,
634 D., and Smith, S. J.: AerChemMIP: quantifying the effects of chemistry and aerosols in CMIP6, *Geosci Model Dev*, 10, 585–607,
635 <https://doi.org/10.5194/gmd-10-585-2017>, 2017.

636 Dagan, G., Stier, P. and Watson-Parris, D.: Contrasting Response of Precipitation to Aerosol Perturbation in the Tropics and
637 Extratropics Explained by Energy Budget Considerations, *Geophys Res Lett*, 46(13), 7828–7837, doi:10.1029/2019gl083479, 2019.

638 Dagan, G., Stier, P. and Watson-Parris, D.: Aerosol Forcing Masks and Delays the Formation of the North Atlantic Warming Hole
639 by Three Decades, *Geophys Res Lett*, 47(22), e2020GL090778, doi:10.1029/2020gl090778, 2020.

640 Danabasoglu, G., Lamarque, J. -F., Bacmeister, J., Bailey, D. A., DuVivier, A. K., Edwards, J., Emmons, L. K., Fasullo, J., Garcia,
641 R., Gettelman, A., Hannay, C., Holland, M. M., Large, W. G., Lauritzen, P. H., Lawrence, D. M., Lenaerts, J. T. M., Lindsay, K.,
642 Lipscomb, W. H., Mills, M. J., Neale, R., Oleson, K. W., Otto-Bliesner, B., Phillips, A. S., Sacks, W., Tilmes, S., Kampenhout, L.,
643 Vertenstein, M., Bertini, A., Dennis, J., Deser, C., Fischer, C., Fox-Kemper, B., Kay, J. E., Kinnison, D., Kushner, P. J., Larson, V.
644 E., Long, M. C., Mickelson, S., Moore, J. K., Nienhouse, E., Polvani, L., Rasch, P. J., and Strand, W. G.: The Community Earth
645 System Model Version 2 (CESM2), *J Adv Model Earth Sy*, 12, <https://doi.org/10.1029/2019ms001916>, 2020.

646 Drijfhout, S., Oldenborgh, G. J. van and Cimadoribus, A.: Is a Decline of AMOC Causing the Warming Hole above the North
647 Atlantic in Observed and Modeled Warming Patterns?, *J Climate*, 25(24), 8373–8379, doi:10.1175/jcli-d-12-00490.1, 2012.

648 Eyring, V., Bony, S., Meehl, G. A., Senior, C. A., Stevens, B., Stouffer, R. J., and Taylor, K. E.: Overview of the Coupled Model
649 Intercomparison Project Phase 6 (CMIP6) experimental design and organization, *Geosci Model Dev*, 9, 1937–1958,
650 <https://doi.org/10.5194/gmd-9-1937-2016>, 2016.

651 Forster, P. *et al.* *The Earth's Energy Budget, Climate Feedbacks, and Climate Sensitivity.* (2021).

652 Gettelman, A., Hannay, C., Bacmeister, J. T., Neale, R. B., Pendergrass, A. G., Danabasoglu, G., Lamarque, J. -F., Fasullo, J. T.,
653 Bailey, D. A., Lawrence, D. M., and Mills, M. J.: High Climate Sensitivity in the Community Earth System Model Version 2
654 (CESM2), *Geophys Res Lett*, 46, 8329–8337, <https://doi.org/10.1029/2019gl083978>, 2019.

655 Gillett, N. P., Shiogama, H., Funke, B., Hegerl, G., Knutti, R., Matthes, K., Santer, B. D., Stone, D., and Tebaldi, C.: The Detection
656 and Attribution Model Intercomparison Project (DAMIP v1.0) contribution to CMIP6, *Geosci Model Dev*, 9, 3685–3697,
657 <https://doi.org/10.5194/gmd-9-3685-2016>, 2016.

658 Gillett, N. P. Kirchmeier-Young, M., Ribes, A., *et al.* Constraining human contributions to observed warming since the pre-
659 industrial period. *Nat Clim Change* 11, 207–212 (2021).

660 Gneiting, T., Raftery, A. E., III, A. H. W. and Goldman, T.: Calibrated Probabilistic Forecasting Using Ensemble Model Output
661 Statistics and Minimum CRPS Estimation, *Mon Weather Rev*, 133(5), 1098–1118, doi:10.1175/mwr2904.1, 2005.

662 Hansen, J., Sato, M., and Ruedy, R.: Long-term changes of the diurnal temperature cycle: implications about mechanisms of
663 global climate change, *Atmos Res*, 37, 175–209, [https://doi.org/10.1016/0169-8095\(94\)00077-q](https://doi.org/10.1016/0169-8095(94)00077-q), 1995.

664 Harder, P., Jones, W., Lguensat, R., Bouabid, S., Fulton, J., Quesada-Chacon, D., Marcolongo, A., Stefanovic, S., Rao, Y.,
665 Manshausen, P. and Watson-Parris, Duncan: NightVision: Generating Night-time Satellite Imagery from Infra-Red Observations,
666 <https://arxiv.org/abs/2011.07017>, 2020

667 Hasselmann, K.: Optimal Fingerprints for the Detection of Time-dependent Climate Change, *J Climate*, 6, 1957–1971,
668 [https://doi.org/10.1175/1520-0442\(1993\)006<1957:offtdo>2.0.co;2](https://doi.org/10.1175/1520-0442(1993)006<1957:offtdo>2.0.co;2), 1993.

669 Ho, T. K.: Random decision forests, *Proc 3rd Int Conf Document Analysis Recognit*, 1, 278–282 vol.1,
670 doi:10.1109/icdar.1995.598994, 1995.

671 Hochreiter, S. and Schmidhuber, J. : Long Short-Term Memory. *Neural computation*, 9(8):1735–1780, 1997.

672 Hoesly, R. M., Smith, S. J., Feng, L., Klimont, Z., Janssens-Maenhout, G., Pitkanen, T., Seibert, J. J., Vu, L., Andres, R. J., Bolt, R.
673 M., Bond, T. C., Dawidowski, L., Kholod, N., Kurokawa, J., Li, M., Liu, L., Lu, Z., Moura, M. C. P., O'Rourke, P. R., and Zhang,
674 Q.: Historical (1750–2014) anthropogenic emissions of reactive gases and aerosols from the Community Emissions Data System
675 (CEDS), *Geosci Model Dev*, 11, 369–408, <https://doi.org/10.5194/gmd-11-369-2018>, 2018.

676 Holden, P. B. and Edwards, N. R.: Dimensionally reduced emulation of an AOGCM for application to integrated assessment
677 modelling: DIMENSIONALLY REDUCED AOGCM EMULATION, *Geophys Res Lett*, 37, n/a-n/a,
678 <https://doi.org/10.1029/2010gl045137>, 2010.

679 Kasoar, M., Shawki, D. and Voulgarakis, A.: Similar spatial patterns of global climate response to aerosols from different regions,
680 *npj Clim Atmospheric Sci*, 1(1), 12, doi:10.1038/s41612-018-0022-z, 2018.

681 Kendall, A. and Gal, Y.: What Uncertainties Do We Need in Bayesian Deep Learning for Computer Vision? In *Proceedings of the*
682 *31st International Conference on Neural Information Processing Systems* (pp. 5580-5590), 2017.

683 Knutti, R., Masson, D., and Gettelman, A.: Climate model genealogy: Generation CMIP5 and how we got there, *Geophys Res Lett*,
684 40, 1194–1199, <https://doi.org/10.1002/grl.50256>, 2013.

685 Jeevanjee, N. and Romps, D. M.: Mean precipitation change from a deepening troposphere., *P Natl Acad Sci Usa*, 115(45), 11465–
686 11470, doi:10.1073/pnas.1720683115, 2018.

687 Le Cun, Y., Boser, B., Denker, J.S., Henderson, D., Howard, R.E., Hubbard, W. , Jackel, L.D. et al. : Hand-written digit
688 recognition with a back-propagation network. In *Advances in neural information processing systems*, 1990.

689 Lee, K., Lee, K., Lee, H., and Shin, J.: A Simple Unified Framework for Detecting Out-of-Distribution Samples and Adversarial
690 Attacks, *Arxiv*, 2018.

691 Li, Z., Kovachki, N., Azizzadenesheli, K., Liu, B., Bhattacharya, K., Stuart, A., and Anandkumar, A.: Fourier Neural Operator for
692 Parametric Partial Differential Equations, *Arxiv*, 2020.

693 Matthews, H. D. and Caldeira, K.: Stabilizing climate requires near-zero emissions, *Geophys Res Lett*, 35(4),
694 doi:10.1029/2007gl032388, 2008.

695 Manabe, S. and Stouffer, R. J.: Century-scale effects of increased atmospheric CO₂ on the ocean–atmosphere system, *Nature*,
696 364(6434), 215–218, doi:10.1038/364215a0, 1993.

697 Mansfield, L. A., Nowack, P. J., Kasoar, M., Everitt, R. G., Collins, W. J., and Voulgarakis, A.: Predicting global patterns of long-
698 term climate change from short-term simulations using machine learning, *npj Clim Atmospheric Sci*, 3, 44,
699 <https://doi.org/10.1038/s41612-020-00148-5>, 2020.

700 McGovern, A., Lagerquist, R., II, D. J. G., Jergensen, G. E., Elmore, K. L., Homeyer, C. R., and Smith, T.: Making the black box
701 more transparent: Understanding the physical implications of machine learning Making the black box more transparent:
702 Understanding the physical implications of machine learning, *B Am Meteorol Soc*, 100, 2175–2199, [https://doi.org/10.1175/bams-d-](https://doi.org/10.1175/bams-d-18-0195.1)
703 18-0195.1, 2019.

704 Meinshausen, M., Raper, S. C. B., and Wigley, T. M. L.: Emulating coupled atmosphere-ocean and carbon cycle models with a
705 simpler model, *MAGICC6 – Part 1: Model description and calibration*, *Atmos Chem Phys*, 11, 1417–1456,
706 <https://doi.org/10.5194/acp-11-1417-2011>, 2011.

707 Millar, R. J., Fuglestedt, J. S., Friedlingstein, P., Rogelj, J., Grubb, M. J., Matthews, H. D., Skeie, R. B., Forster, P. M., Frame, D.
708 J., and Allen, M. R.: Emission budgets and pathways consistent with limiting warming to 1.5 °C, *Nat Geosci*, 10, 741,
709 <https://doi.org/10.1038/ngeo3031>, 2017.

710 Nicholls, Z. R. J., Meinshausen, M., Lewis, J., Gieseke, R., Dommenges, D., Dorheim, K., Fan, C.-S., Fuglestedt, J. S., Gasser, T.,
711 Golüke, U., Goodwin, P., Hartin, C., Hope, A. P., Kriegler, E., Leach, N. J., Marchegiani, D., McBride, L. A., Quilcaille, Y., Rogelj,
712 J., Salawitch, R. J., Samset, B. H., Sandstad, M., Shiklomanov, A. N., Skeie, R. B., Smith, C. J., Smith, S., Tanaka, K., Tsutsui, J.,
713 and Xie, Z.: Reduced Complexity Model Intercomparison Project Phase 1: introduction and evaluation of global-mean
714 temperature response, *Geosci Model Dev*, 13, 5175–5190, <https://doi.org/10.5194/gmd-13-5175-2020>, 2020.

715 O'Neill, B. C., Tebaldi, C., Vuuren, D. P. van, Eyring, V., Friedlingstein, P., Hurtt, G., Knutti, R., Kriegler, E., Lamarque, J.-F.,
716 Lowe, J., Meehl, G. A., Moss, R., Riahi, K., and Sanderson, B. M.: The Scenario Model Intercomparison Project (ScenarioMIP)
717 for CMIP6, *Geosci Model Dev*, 9, 3461–3482, <https://doi.org/10.5194/gmd-9-3461-2016>, 2016.

718 Oreskes, N., Shrader-Frechette, K. and Belitz, K.: Verification, Validation, and Confirmation of Numerical Models in the Earth
719 Sciences, *Science*, 263(5147), 641–646, doi:10.1126/science.263.5147.641, 1994.

720 Otto, F. E. L., Oldenborgh, G. J. van, Eden, J., Stott, P. A., Karoly, D. J., and Allen, M. R.: The attribution question, 6, 813–816,
721 <https://doi.org/10.1038/nclimate3089>, 2016.

722 Pendergrass, A. G. and Hartmann, D. L.: The Atmospheric Energy Constraint on Global-Mean Precipitation Change, *J Climate*,
723 27(2), 130916120136005, doi:10.1175/jcli-d-13-00163.1, 2013.

724 Pincus, R., Batstone, C. P., Hofmann, R. J. P., Taylor, K. E., and Glecker, P. J.: Evaluating the present-day simulation of clouds,
725 precipitation, and radiation in climate models, *J Geophys Res*, 113, <https://doi.org/10.1029/2007jd009334>, 2008.

726 Pincus, R., Forster, P. M., and Stevens, B.: The Radiative Forcing Model Intercomparison Project (RFMIP): experimental
727 protocol for CMIP6, *Geosci Model Dev*, 9, 3447–3460, <https://doi.org/10.5194/gmd-9-3447-2016>, 2016.

728 Rabanser, S., Günemann, S. and Lipton, Z. C.: Failing Loudly: An Empirical Study of Methods for Detecting Dataset Shift,
729 *Arxiv*, 2018.

730 Raissi, M., Perdikaris, P., and Karniadakis, G. E.: Numerical Gaussian Processes for Time-Dependent and Nonlinear Partial
731 Differential Equations, *Siam J Sci Comput*, 40, A172–A198, <https://doi.org/10.1137/17m1120762>, 2018.

732 Rasmussen, C. E. and Williams, C. K. I.: *Gaussian Processes for Machine Learning*, , doi:10.7551/mitpress/3206.001.0001, 2005.

733 Rasp, S., Dueben, P. D., Scher, S., Weyn, J. A., Mouatadid, S., and Thuerey, N.: WeatherBench: A Benchmark Data Set for Data-
734 Driven Weather Forecasting, *J Adv Model Earth Sy*, 12, <https://doi.org/10.1029/2020ms002203>, 2020.

735 Richardson, T. B., Forster, P. M., Smith, C. J., Maycock, A. C., Wood, T., Andrews, T., Boucher, O., Faluvegi, G., Fläschner, D.,
736 Hodnebrog, Ø., Kasoar, M., Kirkevåg, A., Lamarque, J.-F., Mülmenstädt, J., Myhre, G., Olivé, D., Portmann, R. W., Samset, B.
737 H., Shawki, D., Shindell, D., Stier, P., Takemura, T., Voulgarakis, A., and Watson-Parris, D.: Efficacy of Climate Forcings in
738 PDRMIP Models., *J Geophys Res Atmospheres Jgr*, 124, 12824–12844, <https://doi.org/10.1029/2019jd030581>, 2019.

739 Rodgers, K. B., Lee, S.-S., Rosenbloom, N., Timmermann, A., Danabasoglu, G., Deser, C., Edwards, J., Kim, J.-E., Simpson, I. R.,
740 Stein, K., Stuecker, M. F., Yamaguchi, R., Bódai, T., Chung, E.-S., Huang, L., Kim, W. M., Lamarque, J.-F., Lombardozzi, D. L.,
741 Wieder, W. R., and Yeager, S. G.: Ubiquity of human-induced changes in climate variability, *Earth Syst. Dynam.*, 12, 1393–1411,
742 <https://doi.org/10.5194/esd-12-1393-2021>, 2021

743 Ronneberger, O., Fischer, P. and Brox, T. : U-Net: Convolutional Networks for Biomedical Image Segmentation. LNCS. 9351. 234-
744 241. 10.1007/978-3-319-24574-4_28, 2015.

745 Schneider, T., Bischoff, T., and Haug, G. H.: Migrations and dynamics of the intertropical convergence zone, 513,
746 <https://doi.org/10.1038/nature13636>, 2014.

747 Seland, Ø., Bentsen, M., Olivé, D., Toniazzo, T., Gjermundsen, A., Graff, L. S., Debernard, J. B., Gupta, A. K., He, Y.-C.,
748 Kirkevåg, A., Schwinger, J., Tjiputra, J., Aas, K. S., Bethke, I., Fan, Y., Griesfeller, J., Grini, A., Guo, C., Ilicak, M., Karset, I. H.
749 H., Landgren, O., Liakka, J., Moseid, K. O., Nummelin, A., Spensberger, C., Tang, H., Zhang, Z., Heinze, C., Iversen, T., and
750 Schulz, M.: Overview of the Norwegian Earth System Model (NorESM2) and key climate response of CMIP6 DECK, historical,
751 and scenario simulations, *Geosci Model Dev*, 13, 6165–6200, <https://doi.org/10.5194/gmd-13-6165-2020>, 2020.

752 Sellar, A. A., Jones, C. G., Mulcahy, J. P., Tang, Y., Yool, A., Wiltshire, A., O'Connor, F. M., Stringer, M., Hill, R., Palmieri, J.,
753 Woodward, S., Mora, L., Kuhlbrodt, T., Rumbold, S. T., Kelley, D. I., Ellis, R., Johnson, C. E., Walton, J., Abraham, N. L.,
754 Andrews, M. B., Andrews, T., Archibald, A. T., Berthou, S., Burke, E., Blockley, E., Carslaw, K., Dalvi, M., Edwards, J., Folberth,
755 G. A., Gedney, N., Griffiths, P. T., Harper, A. B., Hendry, M. A., Hewitt, A. J., Johnson, B., Jones, A., Jones, C. D., Keeble, J.,
756 Liddicoat, S., Morgenstern, O., Parker, R. J., Predoi, V., Robertson, E., Siahann, A., Smith, R. S., Swaminathan, R., Woodhouse,
757 M. T., Zeng, G., and Zerroukat, M.: UKESM1: Description and Evaluation of the U.K. Earth System Model, *J Adv Model Earth*
758 *Sy*, 11, 4513–4558, <https://doi.org/10.1029/2019ms001739>, 2019.

759 Shindell, D. T., Faluvegi, G., Koch, D. M., Schmidt, G. A., Unger, N., and Bauer, S. E.: Improved Attribution of Climate Forcing
760 to Emissions, *Science*, 326, 716–718, <https://doi.org/10.1126/science.1174760>, 2009.

761 Smith, C. J., Forster, P. M., Allen, M., Leach, N., Millar, R. J., Passerello, G. A., and Regayre, L. A.: FAIR v1.3: a simple
762 emissions-based impulse response and carbon cycle model, *Geosci Model Dev*, 11, 2273–2297, [https://doi.org/10.5194/gmd-11-2273-](https://doi.org/10.5194/gmd-11-2273-2018)
763 2018, 2018.

764 Smith, C. J., Harris, G. R., Palmer, M. D., Bellouin, N., Collins, W., Myhre, G., Schulz, M., Golaz, J. -C., Ringer, M., Storelvmo,
765 T., and Forster, P. M.: Energy Budget Constraints on the Time History of Aerosol Forcing and Climate Sensitivity, *J Geophys Res*
766 *Atmospheres*, 126, <https://doi.org/10.1029/2020jd033622>, 2021.

767 Stott, P. A., Gillett, N. P., Hegerl, G. C., Karoly, D. J., Stone, D. A., Zhang, X., and Zwiers, F.: Detection and attribution of climate
768 change: a regional perspective, *Wiley Interdiscip Rev Clim Change*, 1, 192–211, <https://doi.org/10.1002/wcc.34>, 2010.

769 Stott, P. A., Christidis, N., Otto, F. E. L., Sun, Y., Vanderlinden, J., Oldenborgh, G. J. van, Vautard, R., Storch, H. von, Walton,
770 P., Yiou, P., and Zwiers, F. W.: Attribution of extreme weather and climate-related events, *Wiley Interdiscip Rev Clim Change*, 7,
771 23–41, <https://doi.org/10.1002/wcc.380>, 2016.

772 Tebaldi, C. and Arblaster, J. M.: Pattern scaling: Its strengths and limitations, and an update on the latest model simulations,
773 *Climatic Change*, 122, 459–471, <https://doi.org/10.1007/s10584-013-1032-9>, 2014.

774 Trebing, K., Stanczyk T. and Mehrkanoon, S.: SmaAt-Unet: Precipitation Nowcasting using Small Attention-UNet Architecture,
775 <https://arxiv.org/abs/2007.04417>, 2021.

776 Watson-Parris, D.: Machine learning for weather and climate are worlds apart, *Philosophical Transactions Royal Soc*, 379,
777 20200098, <https://doi.org/10.1098/rsta.2020.0098>, 2021.

778 Watson-Parris, D., Williams, A., Deaconu, L., and Stier, P.: Model calibration using ESEm v1.1.0 – an open, scalable Earth system
779 emulator, *Geosci. Model Dev.*, 14, 7659–7672, <https://doi.org/10.5194/gmd-14-7659-2021>, 2021.

780 Wilcox, L. J., Liu, Z., Samset, B. H., Hawkins, E., Lund, M. T., Nordling, K., Undorf, S., Bollasina, M., Ekman, A. M. L.,
781 Krishnan, S., Merikanto, J., and Turner, A. G.: Accelerated increases in global and Asian summer monsoon precipitation from
782 future aerosol reductions, *Atmos Chem Phys*, 20, 11955–11977, <https://doi.org/10.5194/acp-20-11955-2020>, 2020.

783 Woollings, T., Gregory, J. M., Pinto, J. G., Meyers, M. and Brayshaw, D. J.: Response of the North Atlantic storm track to climate
784 change shaped by ocean–atmosphere coupling, *Nat Geosci*, 5(5), 313–317, doi:10.1038/ngeo1438, 2012.

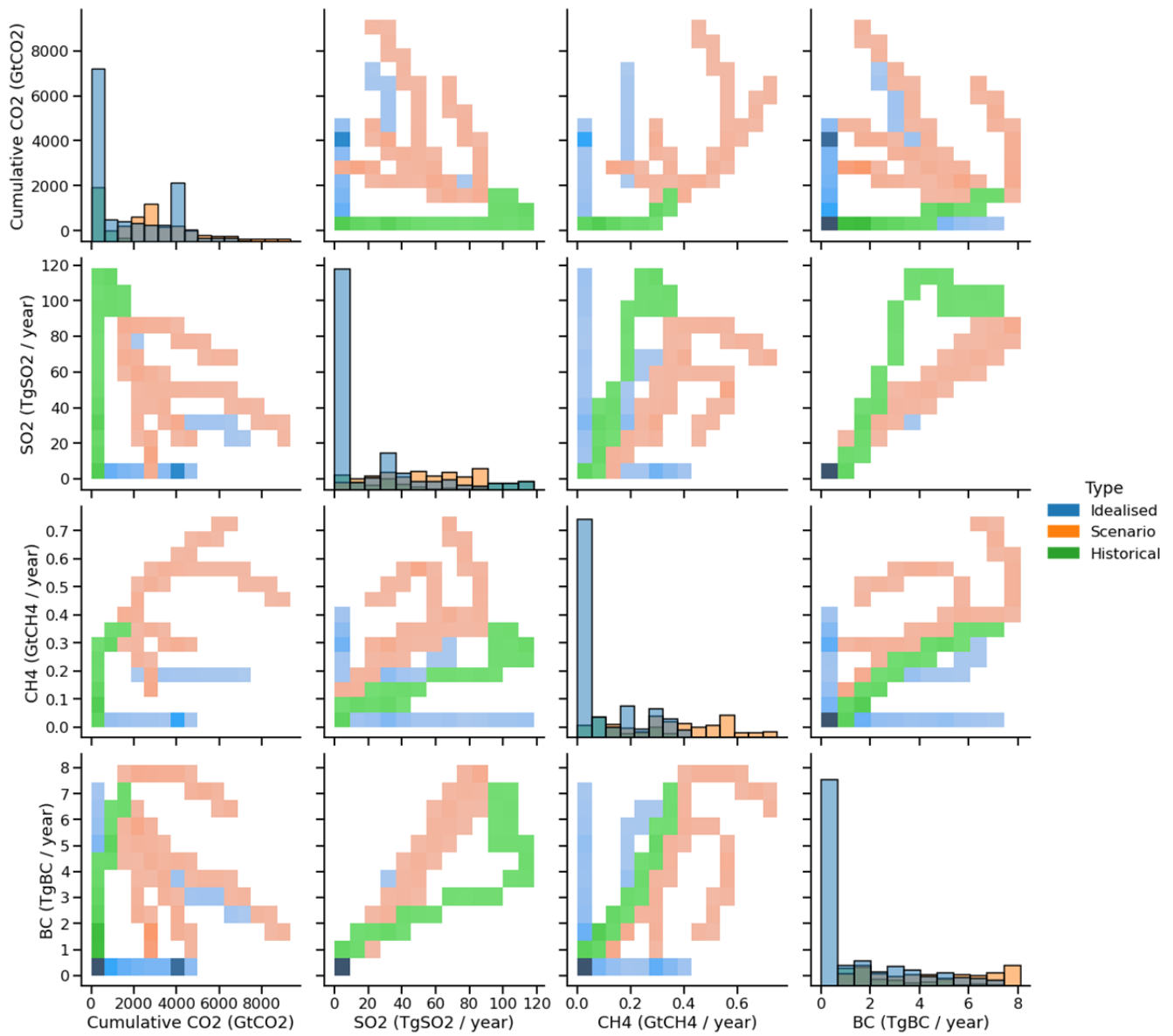
785

786

787

788

789



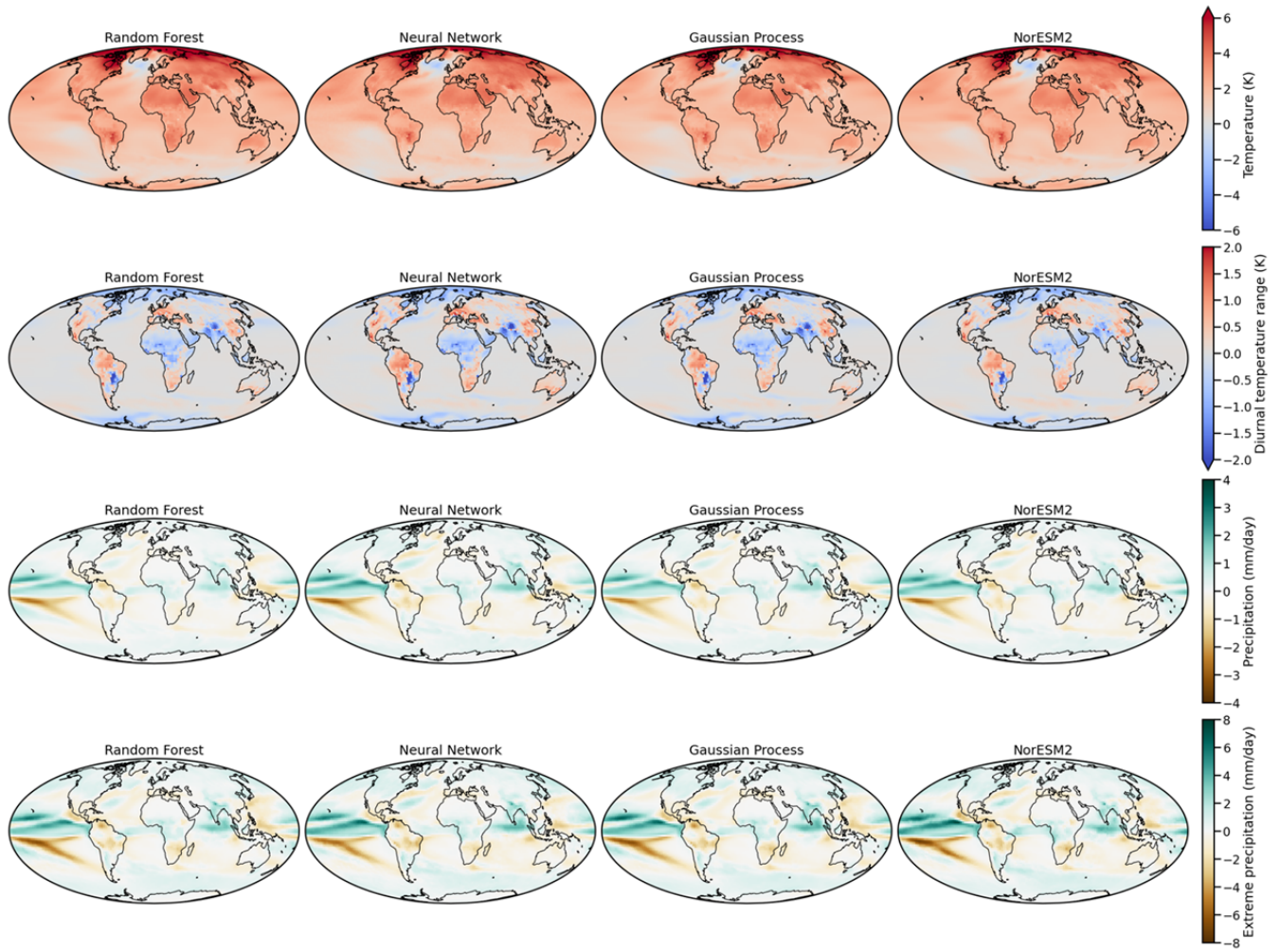
791

792 **Figure A1: Joint and marginal distributions of annual global mean emissions and concentrations across the ClimateBench training**
 793 **dataset. Input datasets are classified as Idealised (such as *IpctCO2* and *abrupt4xCO2*, and including *ssp370-lowNTCF*), Historical**
 794 **and Scenario to demonstrate the contribution of each to sampling the full input space.**

795

796

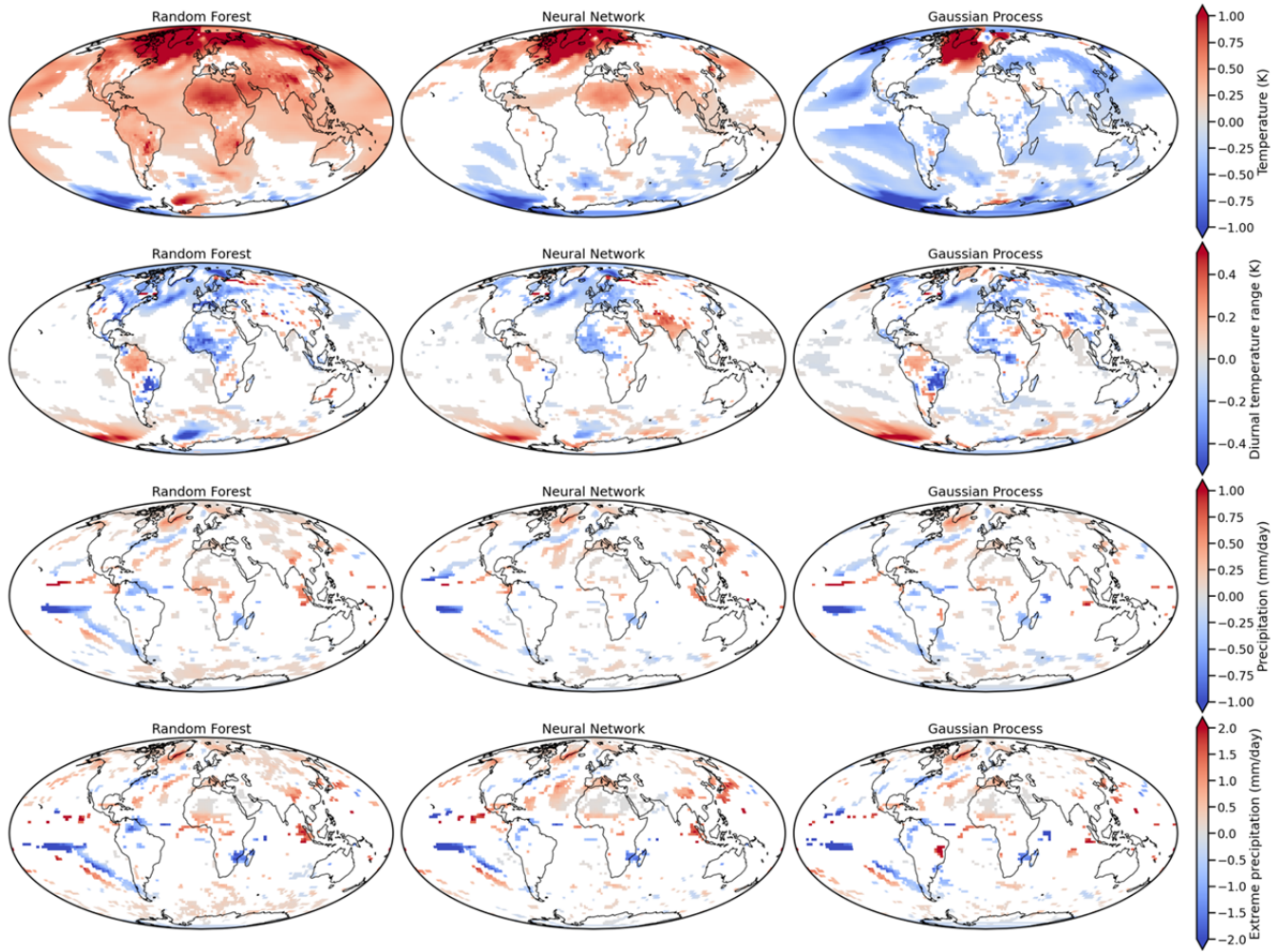
797



798

799 **Figure A2: Maps of ClimateBench target variables for each baseline model and the target NorESM values under the test *ssp245***
 800 **scenario averaged between 2080-2100.**

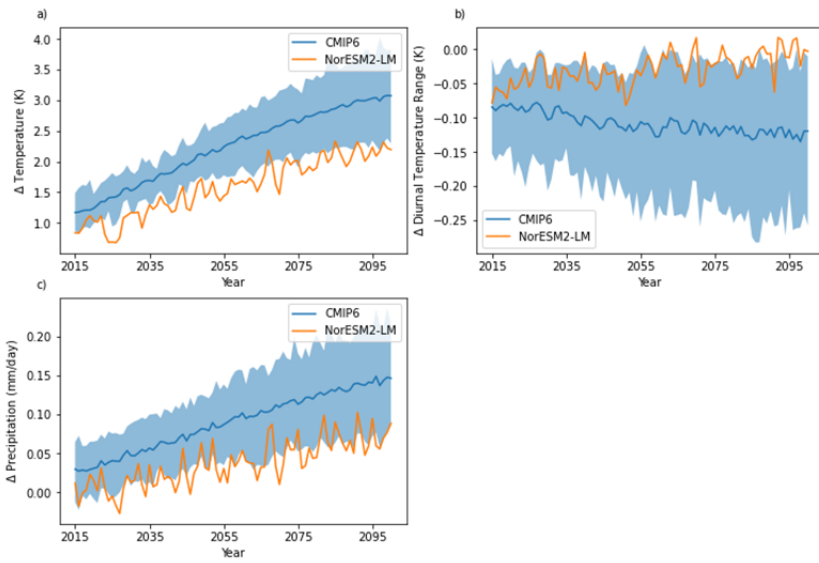
801



802

803 **Figure A3: Maps of the mean difference in the ClimateBench target variables for each baseline emulator against the target**
 804 **NorESM values under the test ssp245-aer scenario averaged between 2080-2100. Differences insignificant at the p<5% level are**
 805 **masked from the plots.**

806



807
808
809

Figure A4: Global mean NorESM-LM projections under *ssp-245* as compared to all other available CMIP6 models for three of the target variables

810

811

812 **Table A4: The spatial, global and total NRMSE of the different baseline emulators for the years 2080-2100 against the**
 813 **ClimateBench task of estimating key climate variables under the idealised future scenario SSP245-aer.**

	NRMSE surface air temperature (1)			NRMSE diurnal temperature range (1)			NRMSE precipitation (1)			NRMSE 90 th percentile precipitation (1)		
	Spatial	Global	Total	Spatial	Global	Total	Spatial	Global	Total	Spatial	Global	Total
Gaussian Process	2.138	1.165	7.963	14.298	2.868	28.636	12.100	0.933	16.767	13.486	1.353	20.252
Neural Network	2.116	1.011	7.173	12.387	2.200	23.386	10.316	0.977	15.199	12.224	1.438	19.414
Random Forest	2.977	2.041	13.182	16.222	3.284	32.642	11.562	1.291	18.017	12.302	1.616	20.382

814

815 **Table A5: The CRPS for the Gaussian process emulator for the years 2080-2100 against the ClimateBench task of estimating key**
 816 **climate variables under future scenario SSP245.**

	CRPS surface air temperature (K)	CRPS diurnal temperature range (K)	CRPS precipitation (mm / day)	CRPS 90 th percentile precipitation (mm / day)
Gaussian Process	0.4765	0.3601	1.0753	1.0029

817

818

819

820 **1. Gaussian process model specifications**

821 The GP models kernel k have the same form for all four climate response variables

822
$$k = k_{CO_2} + k_{CH_4} + k_{BC} + k_{SO_2}$$

823 where k_{CO_2} and k_{CH_4} are kernels that respectively take as inputs CO2 and CH4 emissions. k_{BC} and k_{SO_2} are kernels that take
824 as inputs the 5 principal components of BC and SO2 emission maps respectively, each principal component being rescaled
825 by an independent length scale term. We choose the Matérn-1.5 class of kernel,

826
$$k_X(x, x') = \left(1 + \sqrt{3}d(x, x')\right) \exp\left(-\sqrt{3}d(x, x')\right),$$

827 where X is a general notation for CO2, CH4, BC or SO2, and $d(x, x')$ is a distance between inputs typically given by

828
$$d(x, x') = \sum_i |x_i - x'_i|/l_i.$$

829 l_i is a length scale associated to the i^{th} coordinate x_i . Global CO2 and CH4 emissions are scalar inputs, hence the
830 corresponding distances only involve one length scale parameter. The principal components decompositions of BC and SO2
831 emission maps both have 5 coordinates, hence we set each principal component to be a different coordinate with its own
832 length scale parameter. The Matérn-1.5 kernel guarantees that the corresponding GP lies in a space of continuous functions,
833 hence providing regularity to the climate response predictions. We refer the reader to Rasmussen and Williams, 2005,
834 Chapter 4 for more details on the Matérn kernel. Each kernel is multiplied by a variance term σ_X^2 , which rescales the kernel
835 in the above sum and allows to balance relative features importance. Variances and length scales are tuned during the
836 optimization step.

837

838

2. Random forest model specification

839

Hyperparameter	number of trees	min samples split	min samples leaf	maxdepth
Surface air temperature	250	5	7	5
Diurnal temperature range	150	15	8	40
Precipitation	250	15	12	25
90 th percentile of precipitation	300	10	12	20

840

841

3. Neural Network model specification

842

The parameters are the same for all four models.

843

844

Model architecture

Layer	Hyperparameter value (if not specified, the default parameters are used)	Output Shape	Param #
Time distributed Conv2D	Number of filters: 20 Filter size: 3 Activation function: ReLu	(None, 10, 96, 144, 20)	740
Time distributed AveragePooling2D	Pool size: 2	(None, 10, 48, 72, 20)	0
Time distributed GlobalAveragePooling2D		(None, 10, 20)	0
LSTM	Number of units: 25 Activation function: ReLu	(None, 25)	4600
Dense	Units: 96*144	(None, 13824)	359424
Activation	Activation function: linear	(None, 13824)	0
Reshape		(None, 1, 96, 144)	0

845

846

Model training

Hyperparameter	Value
----------------	-------

Batch size	16
Epochs	30
Optimizer	Rmsprop
Metric	MSE

847

848