Detection of forced change within combined climate fields using explainable neural networks

Jamin Kurtis Rader^{1,1}, Elizabeth A. Barnes^{1,1}, Imme Ebert-Uphoff^{1,1}, and Chuck Anderson^{1,1}

¹Colorado State University

November 30, 2022

Abstract

Assessing forced climate change requires the extraction of the forced signal from the background of climate noise. Traditionally, tools for extracting forced climate change signals have focused on one atmospheric variable at a time, however, using multiple variables can reduce noise and allow for easier detection of the forced response. Following previous work, we train artificial neural networks to predict the year of single- and multi-variable maps from forced climate model simulations. To perform this task, the neural networks learn patterns that allow them to discriminate between maps from different years—that is, the neural networks learn the patterns of the forced signal amidst the shroud of internal variability and climate model disagreement. When presented with combined input fields (multiple seasons, variables, or both), the neural networks are able to detect the signal of forced change earlier than when given single fields alone by utilizing complex, nonlinear relationships between multiple variables and seasons. We use layer-wise relevance propagation, a neural network explainability tool, to identify the multivariate patterns "vary in time and between climate models, providing a template for investigating inter-model differences in the time evolution of the forced response. This work demonstrates how neural networks and their explainability tools can be harnessed to identify patterns of the forced signal within combined fields.

Detection of forced change within combined climate fields using explainable neural networks

Jamin K. Rader¹, Elizabeth A. Barnes¹, Imme Ebert-Uphoff^{2,3}, Chuck Anderson⁴

5	¹ Department of Atmospheric Science, Colorado State University, Fort Collins, CO, USA
6	2 Cooperative Institute for Research in the Atmosphere, Colorado State University, Fort Collins, CO, USA
7	³ Department of Electrical and Computer Engineering, Colorado State University, Fort Collins, CO, USA
8	⁴ Department of Computer Science, Colorado State University, Fort Collins, CO, USA

9 Key Points:

1

2

3

4

10	• Neural networks and their explainability tools can be harnessed to identify pat-
11	terns of forced change within combined fields
12	• Combined fields of input allow for earlier detection of the emergence of a forced
13	climate response
14	• Explainable AI techniques can be used to identify patterns that describe the emer-
15	gence and evolution of forced climate change

 $Corresponding \ author: \ Jamin \ K. \ Rader, \ \texttt{jamin.rader@colostate.edu}$

16 Abstract

Assessing forced climate change requires the extraction of the forced signal from the back-17 ground of climate noise. Traditionally, tools for extracting forced climate change signals 18 have focused on one atmospheric variable at a time, however, using multiple variables 19 can reduce noise and allow for easier detection of the forced response. Following previ-20 ous work, we train artificial neural networks to predict the year of single- and multi-variable 21 maps from forced climate model simulations. To perform this task, the neural networks 22 learn patterns that allow them to discriminate between maps from different years—that 23 is, the neural networks learn the patterns of the forced signal amidst the shroud of in-24 ternal variability and climate model disagreement. When presented with combined in-25 put fields (multiple seasons, variables, or both), the neural networks are able to detect 26 the signal of forced change earlier than when given single fields alone by utilizing com-27 plex, nonlinear relationships between multiple variables and seasons. We use layer-wise 28 relevance propagation, a neural network explainability tool, to identify the multivariate 29 patterns learned by the neural networks that serve as reliable indicators of the forced 30 response. These "indicator patterns" vary in time and between climate models, provid-31 ing a template for investigating inter-model differences in the time evolution of the forced 32 response. This work demonstrates how neural networks and their explainability tools can 33 be harnessed to identify patterns of the forced signal within combined fields. 34

35 Plain Language Summary

Using machine learning tools called neural networks, we identify patterns of the changing climate within climate model data. Changes in the climate can be identified earlier when detecting patterns within maps of multiple variables and seasons than for single maps alone. By visualizing the patterns learned by the neural networks, we can identify which regions, variables, and seasons are most important for detecting climate change. These patterns offer insight into how climate change is represented in different climate models, and how the patterns of climate change will evolve over time.

43 **1** Introduction

Changes in the climate system comprise the Earth system's response to anthropogenic external forcings (e.g. greenhouse gas and aerosol emissions), natural external forcings (e.g. variations in the solar cycle, volcanic activity), internal variability (natural variations in the climate due to internal processes), and the interactions between them. Distinguishing which features of climate change are the product of external forcings, rather than a byproduct of internal variability, is critical for mitigation and adaptation science

(Field et al., 2014; Maher et al., 2021; Mankin et al., 2020; Sanderson et al., 2018). To 50 identify the forced response to external forcings, changes in the climate are often sim-51 plified into "signal" and "noise" components (e.g., Hawkins & Sutton, 2009; Mahony & 52 Cannon, 2018; Scaife & Smith, 2018). The signal of climate change captures all anthro-53 pogenic and natural external forcings, which we refer to as the forced signal or forced 54 response in this study. Climate noise, a combination of internal variability (natural vari-55 ations in the climate due to internal processes) and climate model disagreement in the 56 magnitude of the response, often acts to obscure the forced signal (Santer et al., 2011). 57

Innovative methods are required to determine which behaviors of the climate are 58 the result of the forced signal and which are the result of climate noise. Decades of re-59 search have provided a diverse toolkit for this task (North & Stevens, 1998) which in-60 cludes linear regression (e.g., Mudelsee, 2019; Santer et al., 1996; Sippel et al., 2020; Solow, 61 1987), empirical orthogonal functions and linear discriminant analysis (e.g., Santer et 62 al., 2019; Schneider & Held, 2001; Wills et al., 2018, 2020), and linear inverse models (e.g., 63 Solomon & Newman, 2012), to name a few. Recently, neural networks have also entered 64 the fold. Neural networks are machine learning algorithms that are able to detect com-65 plex, nonlinear relationships between input and output data (Abiodun et al., 2018). Be-66 cause neural networks are able to detect highly complex relationships, they are useful 67 for many high dimensional problems and have become prevalent in several atmospheric 68 science research fields, such as weather forecasting (e.g., Lagerquist et al., 2019; Lee et 69 al., 2021; Weyn et al., 2020), climate model parameterizations (e.g., Brenowitz & Brether-70 ton, 2018; Gettelman et al., 2021; Silva et al., 2021), and, most relevant to the focus of 71 this study, detection of a forced climate response (e.g., Barnes et al., 2019, 2020; Labe 72 & Barnes, 2021; Madakumbura et al., 2021). To detect patterns of forced change, Barnes 73 et al. (2020) trained a neural network to predict the year label of maps of annual-mean 74 temperature (or precipitation) from climate model simulations for forced historical and 75 future scenarios. Given that the internal variability in any given year differs between the 76 various climate models, the neural network had to learn patterns of the forced climate 77 response. Using neural network explainability methods, they then visualized the regions 78 that were most reliable indicators for identifying change across the CMIP5 models. Barnes 79 et al. (2020) demonstrated that neural networks, and their explainability methods, are 80 powerful tools for extracting forced patterns from climate data. This neural network method 81 is a natural approach for isolating the forced climate response. While many other meth-82 ods require assumptions to be made about the time evolution of the forced signal and 83 internal variability within the system, neural networks do not (Barnes et al., 2019). Fol-84 lowing Barnes et al. (2020), neural networks have since been used to explore the sensi-85 tivity of regional temperature signals to aerosols and greenhouse gases using single-forcing 86

large ensembles, and to detect the signal of extreme precipitation in observational datasets
(Labe & Barnes, 2021; Madakumbura et al., 2021).

Though many climate signal detection studies focus on single variables, such as annual-89 mean temperature or a single season of precipitation (Gaetani et al., 2020; Li et al., 2017; 90 Santer et al., 1996, 2019), there are benefits to studying climate change through a mul-91 tivariate lens (Bindoff et al., 2013; Bonfils et al., 2020; Mahony & Cannon, 2018). Many 92 variables in our atmosphere are closely interconnected, so when the variables are intel-93 ligently selected signals of change within multiple variables may be detected earlier than 94 in single variables alone. For example, departure from natural variability can be seen decades 95 earlier in bivariate maps of summertime temperature and precipitation than in either 96 variable alone (Mahony & Cannon, 2018). Similarly, Fischer and Knutti (2012) found 97 that climate model biases in the signal of relative humidity and temperature are neg-98 atively correlated such that climate model simulations of their combined quantity, heat 99 stress, have considerably less spread. Combined variables have also been used to iden-100 tify the impacts of anthropogenic forcings on climate in observational datasets by iden-101 tifying the multivariate patterns that enhance the signal of change relative to the un-102 derlying noise (e.g., Barnett et al., 2008; Marvel & Bonfils, 2013). Understanding how 103 the patterns of the forced response take shape through multiple atmospheric variables 104 also allows for a deeper understanding of the physics at play, as in Bonfils et al. (2020). 105 They explored the evolution of the climate fingerprint by analyzing the leading combined 106 empirical orthogonal functions of temperature, precipitation, and climate moisture in-107 dex. This multivariate approach illuminated two cross-variable patterns of change: in-108 tensification of wet-dry patterns and meridional shifts in the ITCZ associated with in-109 terhemispheric temperature contrasts. Neither pattern can be fully explained by a sin-110 gle variable which highlights the utility of combining variables when identifying patterns 111 of the forced response. 112

Combining fields can be useful for identifying patterns of forced change that do not 113 reveal themselves in single fields alone, but this added information does not come with-114 out its drawbacks. Many variables covary in complex and nonlinear ways, such as sea 115 surface temperature and precipitation (Lu et al., 2015), drought indices (Wu et al., 2017), 116 and snowpack, soil moisture and flood risk (Swain et al., 2020), often requiring complex 117 statistics to isolate these interactions. Identifying nonlinear correlations within climate 118 fields introduces another issue, namely in explaining the complex interplay between fields. 119 These drawbacks highlight the need for methods that are both complex and explainable 120 in multivariate climate analyses. 121

Providing a method for both nonlinear and multi-variable analysis of the forced re-122 sponse, this study extends the neural-network approach of Barnes et al. (2020) to com-123 bined fields of input. Combined fields could mean the same variable for different tem-124 poral segments (e.g. seasons), or different geophysical variables, both of which are ex-125 plored here. For the sake of consistency and comparability, this study largely follows the 126 methodology of Barnes et al. (2020), however there are some departures. We standard-127 ize the input fields differently which improves the predictive skill of the neural networks. 128 We also use a slightly simpler neural network architecture to reduce the computational 129 expense of training a single neural network, and the results from multiple neural networks, 130 rather than just one, are explored. Barnes et al. (2020) demonstrated the utility of neu-131 ral network explainability methods, and we use these methods in tandem with a clus-132 tering technique to enhance post-hoc explanations of neural network decisions. 133

Section 2 outlines the climate models and observations analyzed in this study. Sec-134 tion 3 introduces the neural network design, the explainability technique (layer-wise rel-135 evance propagation; LRP), and their applications to detection of the forced climate re-136 sponse. We then apply these methods to global temperature and precipitation over land 137 in Section 4. Here we investigate the benefits of combining variables and compare the 138 results of the neural network with the classical approach of calculating signal-to-noise 139 ratios. In Section 5, we explore the patterns of the forced response for extreme precip-140 itation over the Americas and investigate the applications of LRP to studying the evo-141 lution of nonlinear climate patterns across multiple climate models. Finally, Section 6 142 summarizes the results of this work and its implications for future work in forced change 143 detection. 144

145 **2 Data**

146

2.1 CMIP6 Climate Models

We use climate model output from the sixth phase of the Coupled Model Intercom-147 parison Project (CMIP6; Eyring et al., 2016). Specifically we focus on monthly-, seasonal-148 , and annual-mean fields of 2-meter air temperature (K), precipitation rate $(kg m^{-2} s^{-1})$, 149 and precipitation rate from very wet days $(kg m^{-2} s^{-1})$, hereafter referred to as temper-150 ature, precipitation, and extreme precipitation, respectively. We use the meteorological 151 seasons of December-January-February (DJF), March-April-May (MAM), June-July-August 152 (JJA), and September-October-November (SON) for calculating seasonal-mean fields. 153 Defining seasons in this way allows for the earliest detection of forced change (see Fig-154 ure S1 for more details). 155

-5-

manuscript submitted to Journal of Advances in Modeling Earth Systems (JAMES)

Very wet days are defined as days that exceed the 95th percentile of all days with 156 precipitation over a pre-defined baseline period (Donat et al., 2016). This is a popular 157 index for measuring changes in extreme precipitation (Cui et al., 2019; Kim et al., 2020) 158 and is used as an indicator of climate change in the U.S. Global Climate Research Pro-159 gram (USGCRP, 2018). We define the baseline as the 40 years from 1980 to 2019, a pe-160 riod for which daily precipitation data exists in both the climate models and the obser-161 vations. To remove the instances in which climate models simulate sub-trace daily pre-162 cipitation totals, we only include days that simulated at least 1 mm of precipitation when 163 calculating the 95th percentile of all days with precipitation (Dai et al., 2007). 164

The neural networks are trained on CMIP6 climate model data. One ensemble mem-165 ber is selected for each of the 37 CMIP6 climate models analyzed so each climate model 166 is only represented once in the training and testing data. Since daily output is required 167 to calculate very wet days, we are limited to 32 models for extreme precipitation (Fig-168 ure S3). We analyze the climate model data from 1920 to 2098 under historical forcing 169 (1920–2014) and the SSP585 scenario (2015–2098). SSP585 represents the highest de-170 velopment pathway within CMIP6 scenarios (O'Neill et al., 2016), combining shared so-171 cioeconomic pathway 5 (SSP5) and representative concentration pathway 8.5 (RCP8.5). 172

Our neural network methodology requires that all climate model fields have the same 173 shape. To accommodate this we regrid the climate model fields from their native res-174 olutions using the second-order conservative remapping method in the Climate Data Op-175 erators package from MPI (Schulzweida, 2019). This regridding step reduces the spatial 176 resolution of the data for most climate models. For temperature and precipitation, the 177 data is regridded to 4 degrees latitude by 4 degrees longitude. We elect to use lower res-178 olution data to reduce the computational expense of training neural networks over global 179 maps of temperature and precipitation. Since the domain for extreme precipitation is 180 smaller than the domain for temperature and precipitation (see the following paragraph), 181 and higher resolution data may better capture regional extreme precipitation patterns, 182 the data for extreme precipitation is regridded to a slightly higher resolution: 1.5 degrees 183 latitude by 1.5 degrees longitude. 184

Two spatial domains are considered in the results of this paper. For temperature and precipitation, the neural networks are trained on all land north of 60°S. Here, we choose to focus on land grid points because that is where humanity lives and will acutely feel the impacts of changing surface temperatures and precipitation. We also exclude Antarctica where climate models and reanalyses struggle to accurately simulate temperature and precipitation. Each map of temperature and precipitation has 948 unique data points. For extreme precipitation, the neural networks are trained on North and South Amer-

-6-

ica (land grid points bounded by 90°N, 55°S, 170°W, and 25°W). Here, we choose to narrow the regional scope to show that neural networks are powerful tools for identifying
the forced response even when the spatial domain, and thus the available data, is limited. Each map of extreme precipitation has 2314 unique data points.

2.2 Observations

While this work largely focuses on the results of neural networks trained and tested 197 on climate model data, we show that neural networks trained on climate model data can 198 be applied to observational data as well. For temperature, we use the Berkeley Earth 199 Surface Temperature (BEST) dataset (Rohde & Hausfather, 2020). This dataset pro-200 vides both a temperature climatology and the anomalies at monthly resolution from 1850 201 to the present. We added the anomalies to the climatology to reconstruct the absolute 202 temperature (K) at each grid point for all months between 1920 and 2019. Monthly ob-203 servational precipitation fields are obtained from the NOAA Global Precipitation Cli-204 matology Project (GPCP), version 2.3, for 1979 to the present (Adler et al., 2018). Since 205 daily precipitation fields are required to calculate extreme precipitation, and daily GPCP 206 precipitation observations are only available back to October 1996, we elected to calcu-207 late observed extreme precipitation using the European Centre for Medium-Range Weather 208 Forecasts' ERA5 global reanalysis (Hersbach et al., 2020) at 6-hour resolution from 1980 209 to present. All observations are regridded in the same way as the climate model data for 210 each respective variable. 211

3 Forced Change Detection Framework

213

196

3.1 Neural Network Design

To identify indicator patterns of the forced response for combined fields we first de-214 velop artificial neural networks that, given maps of CMIP6 climate model output from 215 every simulated year from 1920 to 2098, are tasked to predict the year that is being sim-216 ulated. The results for neural networks trained on ten different input vectors are explored 217 in the following two sections. The input vectors include annual-, seasonal-, and monthly-218 mean data for temperature, precipitation, and temperature and precipitation combined, 219 as well as seasonal-mean maps for extreme precipitation over the Americas. We use this 220 diverse selection of input vectors to compare neural network performance and indicator 221 patterns for single-field and combined-field inputs. 222

The neural network architecture is illustrated in Figure 1. Each unit of the input layer corresponds to a different grid point in the input fields. For example, a neural net-

manuscript submitted to Journal of Advances in Modeling Earth Systems (JAMES)

work that uses seasonal-mean maps of temperature and precipitation as input (two variables and four seasons for a total of eight maps, 948 grid points per map) would have an input vector with 7,584 units. In all cases, this input layer is followed by two fully connected hidden layers with ten nodes each. The hidden layers are followed by an output layer that consists of 22 classes, one corresponding to each decade midpoint between 1905 and 2115 (e.g. 1905, 1915, 1925, ..., 2115). A softmax function is applied to the outputs to convert them to units of likelihood, where the sum of the output vector is one.

This is a relatively simple architecture for a neural network. Neural networks with 232 this architecture learn the patterns of forced change well, and more complicated archi-233 tectures do not substantially improve neural network performance (see Figure S2). It is 234 also notable that this neural network architecture performs better than multiple linear 235 regression, especially when trained on precipitation, and thus using nonlinear techniques 236 improves our ability to detect the year via patterns of forced change (Figure S2). This 237 simple architecture is also widely accessible to most in the climate science community 238 as it can be trained on a personal laptop-highly complex architectures can be prohibitively 239 computationally expensive (Chen et al., 2020). These neural networks were trained on 240 a standard desktop computer with 16 GB of RAM and a 3.1 GHz, 6-core processor. Train-241 ing a single network took anywhere between two and ten minutes depending on the size 242 of the input field. More details on the neural network design and hyperparameter tun-243 ing can be found in the supplementary materials. 244

The neural network is tasked with "predicting the year" rather than "predicting 245 the decade" as the output layer may suggest. To translate between decade midpoints and 246 individual year labels, we use fuzzy encoding (Zadeh, 1965) such that each year can be 247 mapped to one or more neighboring classes with varying degrees of membership (encoded 248 as likelihood). This is different than traditional methods that would map each year to 249 a single decade midpoint. In the traditional case, 2040 and 2049 would be considered to 250 be members of the same class since they are in the same decade, and information would 251 be lost as there is no way to distinguish whether the samples come from the beginning 252 or the end of the decade. Using fuzzy encoding, this information of where a sample lies 253 in each decade is retained. We use a triangular membership function (Zadeh, 1965) with 254 a width equal to one decade such that each year has partial membership in one or two 255 neighboring decade classes, and the total membership sums to one. Following this method, 256 any year directly on a decade midpoint has membership in that class only while years 257 that fall between decade midpoints have membership in the two neighboring classes. The 258 year 1925, for example, is mapped to a likelihood of one for the class 1925 and a like-259 lihood of zero in all other classes. The year 2078 is mapped to a likelihood of 0.7 for the 260

-8-

261 2075 class and a likelihood of 0.3 for the 2085 class. Note that decoding class likelihoods 262 back to their year is simply the decade-weighted sum of the likelihood: $0.7 \times 2075 + 0.3$ 263 $\times 2085 = 2078$. A visualization of the encoding/decoding process can be found in Fig. 264 2 of Barnes et al. (2020).

265

3.2 Neural Network Training

For each input vector we train 100 neural networks that differ only in which cli-266 mate models are randomly split into the training and testing sets. Partitioning so that 267 each climate model's samples are all part of either the training set or the testing set avoids 268 issues with autocorrelation (i.e. near-identical data appearing in both the training and 269 testing sets). One hundred neural networks provide a range of results across multiple com-270 binations of training and testing simulations and offer confidence that the results are con-271 sistent across CMIP6 climate models and do not overfit to any one training set. Each 272 neural network is trained over the entire 1920-2098 period on 80% of the climate model 273 simulations, and then tested on the remaining 20%. This leads to a training set of 30274 simulations and a testing set of 7 simulations for temperature and precipitation fields, 275 and a training set of 26 simulations and a testing set of 6 simulations for extreme pre-276 cipitation fields. We train the neural networks using the binary cross-entropy loss (see 277 Barnes et al., 2020) between the predicted class likelihoods and the correct class mem-278 bership weights, such that the loss function is minimized when the two are equal. Prop-279 erties of the neural network training process, such as the learning rate and activation func-280 tions, can be found in the supplementary materials. 281

The neural networks have several hidden nodes which enable them to learn com-282 plicated relationships between the input and output data. However, with limited train-283 ing data, many of these learned relationships will capture patterns of the noise in the 284 training dataset which can lead to overfitting (Srivastava et al., 2014). Atmospheric sci-285 ence data is also highly correlated in space and this collinearity can cause complications 286 in the interpretation of the learned weights (Newell & Lee, 1981). Thus, to reduce over-287 fitting and address these issues, we apply ridge regularization (L_2 regularization, see Barnes 288 et al., 2020) to the weights of the first hidden layer. Ridge regularization adds a penalty 289 (called the ridge penalty) to the square of the weights so the solution is penalized for hav-290 ing large weights. Through training, this acts to shrink the largest weights, thus spread-291 ing the weight out more evenly across multiple grid points. In our application this re-292 sults in a more even distribution of importance across regions with strong spatial cor-293 relation and improves the performance of the neural networks when given data they were 294

-9-

not trained on, namely those models in the testing set (elaborated on in Fig. 3, Section
4 of Barnes et al., 2020).

Unlike classical approaches which tune the neural network to reduce the mean squared 297 error (MSE) between the predicted and truth outputs in the testing set (in our case this 298 would be the MSE between the truth and predicted years), we select the ridge penalty 299 that minimizes the time of emergence of the forced climate signal (see Section 3.3). Us-300 ing time of emergence, rather than MSE, to identify the appropriate ridge penalty en-301 sures that we are encouraging the neural networks to learn the patterns of the forced re-302 sponse across all decades. When a small ridge penalty is used, the neural networks are 303 able to predict the year at the end of the 21st century almost perfectly, at the expense 304 of the predictive skill in earlier decades. This results in a later calculation of time of emer-305 gence for the testing set. Slightly increasing the ridge penalty can allow the neural net-306 works to detect the climate change signal slightly earlier (Figure S4). The ridge penalty 307 used for each input vector can be found in the supplementary materials. We use the same 308 ridge penalty for all 100 neural networks trained on each input vector. 309

All input fields (for climate models and observations) are standardized to assist with 310 the training and overall performance of the neural network. We subtracted the 1980–2019 311 mean at each grid point of the input fields for each climate model independently. This 312 recasts each input field to measure the change relative to the 1980–2019 mean, rather 313 than the raw magnitudes, which improves the predictive skill of the neural networks and 314 is also appropriate for identifying indicator patterns of forced change. Since values for 315 precipitation change are often on the order of 10^{-6} , while the values for temperature change 316 are on the order of 10^{0} , we normalized the data so the inputs to the neural network all 317 have a similar magnitude. To do this, the data from 1980–2019 at each grid point for 318 each climate model are detrended using ordinary least squares linear regression. We then 319 take the multi-model mean of the standard deviation of the detrended 1980–2019 data 320 for each grid point. The input fields are then divided by this new field of standard de-321 viations so the inputs are of the same magnitude and fall in a reasonable range for train-322 ing the neural networks. Since all our observational datasets include the years 1980 to 323 2019, we standardize the observations as if they were additional climate models: raw ob-324 servations are subtracted by their own 1980–2019 mean, and divided by the same multi-325 model standard deviations that were used to standardize the CMIP6 data. 326

327

3.3 Time of Emergence Calculation

The time of emergence of the forced climate response (hereafter, simply "TOE") is the time in which the forced response signal is distinguishable from the background

-10-

climate by the neural network. Specifically, we define the TOE as the year when the neu-330 ral network is able to distinguish that year's map from any map over a historical base-331 line period. In this work we define this baseline period as 1920–1959 and, under this def-332 inition, the earliest possible TOE estimate is 1960. The TOE is estimated for each cli-333 mate model simulation independently and a schematic of how the TOE is estimated is 334 presented in Figure 2. First, we calculate the maximum of the neural network-predicted 335 vears over 1920–1959 for each model, which is referred to as the baseline maximum. We 336 then identify the TOE as the earliest year in which a map, and all subsequent maps, per-337 manently exceed the baseline maximum. In Figure 2, sample model 1 has a baseline max-338 imum of 1966 and permanently exceeds this prediction threshold in 2028. Sample model 339 2 has a baseline maximum of 1981 and permanently exceeds this threshold in 1989. Thus, 340 the TOE for sample models 1 and 2 are estimated as 2028 and 1989, respectively. In the 341 following sections we present the TOE for the testing set, however TOE estimates are 342 similar for both the training and testing sets. 343

344

3.4 Layer-wise Relevance Propagation

To visualize the patterns learned by the neural network we apply layer-wise relevance propagation (LRP) which highlights the regions that were most relevant in the neural network's decision-making process (Bach et al., 2015; Montavon et al., 2019). Toms et al. (2020) discusses in detail how LRP can be used for neural network explainability in the geosciences, though the most relevant details of LRP are described here.

LRP is a neural network explainability method that traces how information flows 350 through the pathways of a trained neural network. The values in a single-sample input 351 vector (in our case, a single year) are passed forward through the neural network. Us-352 ing the same weights and activations used in the forward pass, LRP then propagates a 353 single-valued output back through the neural network to infer the extent to which each 354 of the values in the input layer contribute to the output (see Fig. 2 in Bach et al., 2015). 355 We refer to this quantity as relevance. Through this backpropagation process the out-356 put value is conserved such that the sum of all relevance is equal to the output. At first 357 order, relevance can be likened to the product of the regression weights and input map 358 in a linear model. This quantity is natively unitless, but we convert it to a fraction by 359 dividing by the output value. This way, we can consider the relevance of a single pixel 360 in terms of its fractional contribution to the predicted class. Since LRP propagates only 361 a single output value at a time, we propagate relevance only for the decade class with 362 the highest likelihood. While the relevance maps detected by these networks evolve from 363

year to year, this evolution is slow so we find visualizing the highest likelihood decadeis sufficient.

There are several LRP decomposition rules which provide different methods of vi-366 sualizing neural networks (Lapuschkin, 2019; Mamalakis et al., 2021). In our applica-367 tions we use the $\alpha\beta$ -rule which propagates positive relevance (regions that act to increase 368 the class likelihood) and negative relevance (regions that act to decrease the class like-369 lihood) separately. Using the parameters $\alpha = 1$ and $\beta = 0$ we choose to only propa-370 gate positive relevance, thus highlighting the regions that added to the likelihood of the 371 selected decade class. We also looked at the relevance maps for $\beta = 1$ and found that 372 propagating negative relevance did not impact the conclusions. 373

374

3.5 Signal-to-Noise Ratio Calculation

In Section 4, we compare the LRP relevance maps to maps of signal-to-noise ra-375 tio (S/N ratio), a more conventional method for identifying indicator patterns of the forced 376 response. S/N ratio consists of three distinct components: the forced signal, which is di-377 vided by the sum of noise due to internal variability, and noise due to climate model dis-378 agreement. A higher S/N ratio indicates that the signal of the forced response within 379 the climate models is very large relative to the underlying noise. We evaluate the S/N 380 ratio for each grid point separately, following the methodology in Hawkins and Sutton 381 (2012). First, we smooth the data from 1920 to 2098 for each climate model using a fourth-382 order polynomial fit. The signal is defined as the difference between 2090 and 1920 in 383 the smoothed data, while internal variability is defined as the standard deviation of the 384 residuals from the smoothed data, and climate model disagreement is defined as the stan-385 dard deviation of the signals calculated for all the climate models. S/N ratio is calcu-386 lated by dividing the climate signal by the 90% confidence interval in the noise: inter-387 nal variability and climate model disagreement. S/N ratio, and its components, can be 388 seen in Figure S8. 389

- ³⁹⁰ 4 Global Precipitation and Temperature
- 391

4.1 Time of Emergence

Across all input vectors of temperature and precipitation, the neural networks are able to learn patterns of the forced response. In the early-to-mid 20th century the forced signal is small and undetectable by the neural networks amidst the noise of internal variability and model disagreement, which leads to poor predictive skill (Figure 3). However, as the signal increases in magnitude into the late-20th and 21st centuries, the neu-

-12-

ral networks are able to detect the patterns of the forced response and distinguish be-397 tween maps in different years. These patterns of the forced response detected by the neu-398 ral networks are generalizable across CMIP6 models, and as a result the neural network 399 has predictive skill for seen data (the training set, see the supplementary materials) as 400 well as unseen data (the testing set). These behaviors are shown in Figure 3 which presents 401 the predicted years from one trained neural network for each combination of global pre-402 cipitation and temperature input fields. Across all input vectors, a similar story of the 403 forced signal unfolds. Prior to the TOE, the neural network is unable to identify pat-404 terns that allow it to accurately predict the year. As a result, the neural network is equally 405 confident (or unconfident) about which year, between 1920 and the TOE, each input came 406 from, so it predicts years right around the middle of the 20th century. After the TOE, 407 the predicted years tend to follow a 1:1 line with the truth years, indicating that the neu-408 ral network has identified reliable indicators of change for this period. 409

Although the neural networks are trained on climate model simulations, their learned 410 patterns can be used to predict the year for observational data as well. When observa-411 tions are used as input, the predicted years increase with time, just as they do for cli-412 mate model input (Figure 3). This means that the indicators of change derived by the 413 neural networks trained on climate models simulations are largely consistent with the 414 real world. Pearson correlations (r) of the actual years with the years predicted by each 415 neural network are shown in Figure 4. All correlations are positive, indicating that the 416 years predicted by the neural networks increase with time. These correlations are strongest 417 for temperature and combined observations (r ≈ 0.9), but still quite high for precipita-418 tion (r ≈ 0.8). Correlations of actual years with predicted years are slightly higher for 419 the combined temperature and precipitation observations than for temperature obser-420 vations alone (Figure S5), suggesting that the multivariate indicator patterns derived from 421 climate model data are useful for understanding trends in the present-day climate. Across 422 all variables, the highest observational correlations are found by the neural networks trained 423 on seasonal-mean data. The correlation of actual years with predicted years for precip-424 itation observations are sensitive to the dataset of choice, which is expanded on in Sec-425 tion S4 and Figures S5 and S6. 426

The average TOEs, calculated from the climate models in the testing sets of all 100 trained neural networks for each input field (Figure 5), reveal that the forced response can be detected earlier in maps of temperature than in maps of precipitation (Figure 5ac). When presented with combined fields the neural networks are, in many cases, able to detect the forced signal even earlier than when given single fields alone (Figure 5b,f). The TOE is generally earlier for the neural networks trained on seasonal-mean data than

-13-

for the neural networks trained on annual-mean data (Figure 5d-f). This is most notable 433 for precipitation fields, likely because there are large seasonal precipitation responses muted 434 by taking the annual mean (Tabari & Willems, 2018; Zappa et al., 2015). The TOEs are 435 earlier for temperature and precipitation combined than temperature alone when using 436 seasonal-mean maps (Figure 5b), but are approximately equal when using annual-mean 437 or monthly-mean maps (Figure 5a,c), which suggests that precipitation only improves 438 upon the detectability of the forced temperature signal when seasonal-mean fields are 439 used. While annual-mean precipitation may mute seasonal precipitation signals, monthly-440 mean precipitation is noisy. In this case, seasonal means emerge as the appropriate tem-441 poral segments for detecting precipitation change, underlining the importance for the 442 intentional and intelligent selection of neural network inputs. 443

The neural networks identify the earliest TOEs when trained on seasonal-mean tem-444 perature and precipitation combined (Figure 5b,f). The TOE results for all 100 seasonal-445 mean neural networks are summarized in the box plots in Figure S7. While the improve-446 ment in forced response detection is small when precipitation is combined with temper-447 ature, it is still notable given that the forced signal of temperature is much clearer than 448 the forced signal of precipitation. We use these variables as an initial example for em-449 ploying this neural network methodology. We anticipate that more robust results might 450 be found for combinations of variables that have more distinct combined signals, such 451 as humidity and temperature (Fischer & Knutti, 2012). 452

453

4.2 Indicator Patterns for Combined Variables

Having shown that the neural networks are able to predict the year given seasonal 454 means of temperature and precipitation (Figures 3, 5), we now identify and explore the 455 spatial indicator patterns used by the neural networks to make correct predictions. By 456 understanding the neural networks' decision-making process, we can identify which re-457 gions act as combined (multi-seasonal and multi-variable) indicators of forced change amidst 458 a background of internal variability and climate model disagreement. To identify these 459 indicator patterns, we apply LRP to all climate model samples in the training and test-460 ing sets from the year 2090 for the seasonal-mean combined neural networks. Averag-461 ing the LRP results for each season and variable, we highlight the regions that have the 462 highest mean relevance across the 37 CMIP6 climate models and 100 trained neural net-463 works. The relevance maps for temperature (precipitation) are shown in Figure 6a-d (7a-464 d) and indicate the importance of each region in the neural networks' predictions of the 465 year 2090. 466

LRP identifies temperature over North Africa and Central Asia in JJA (Figure 6c) 467 and the Andes and Central Africa in SON (Figure 6d) as the most relevant regions for 468 predicting the year. For precipitation, the regions of highest relevance can be found in 469 Canada and Russia in DJF and SON (Figure 7a,d) and in Central Africa and India in 470 JJA and SON (Figure 7c,d). That is to say that these are the regional patterns iden-471 tified by the neural networks that indicate the presence of forced change across the CMIP6 472 climate models. The scale of the color bars are different between Figures 6 and 7, such 473 that the darkest regions in the temperature maps are approximately one order of mag-474 nitude more relevant than the darkest regions in the precipitation maps. Hence, the neu-475 ral network is relying more heavily on the temperature inputs than the precipitation in-476 puts in order to accurately predict the year. This is not surprising because the forced 477 signal of temperature is clearer than the forced signal of precipitation (Fig. SPM.7 in 478 Field et al., 2014). Even so, including seasonal precipitation allows the neural networks 479 to detect forced change earlier within combined fields than in temperature fields alone 480 (Figure 5b). The improvement in neural network performance provided by precipitation 481 (alongside temperature) is particularly noteworthy given that the S/N ratio for temper-482 ature is larger than the S/N ratio for precipitation in all seasons and regions (Figures 483 6e-h, 7e-h, discussed further in this section). In other words, the forced temperature sig-484 nal is always more pronounced than the forced precipitation signal, but the precipita-485 tion signal is still useful for detecting forced change. 486

LRP is designed to highlight the regions that were most relevant for predicting the 487 correct class (in our case, the correct decade class). These LRP indicator patterns for 488 2090 are not the time-mean patterns of the forced response, they are the patterns used 489 by the neural network to distinguish the end of the 21st century from all other decades. 490 This is distinctly different from S/N ratio which identifies the regions where the forced 491 change from 1920 to 2090 is largest relative to internal variability and climate model spread. 492 Maps of S/N ratio for temperature are shown in Figure 6e-h, and the corresponding maps 493 for precipitation are shown in Figure 7e-h, where a higher S/N ratio (darker green) in-494 dicates a clearer forced signal. These regions of high S/N ratio are consistent with other 495 related studies (e.g., Hawkins et al., 2020). For the most part, the indicator patterns iden-496 tified by LRP correspond with the regions with the highest S/N ratios. Calculating the 497 Spearman's rank correlation (ρ) between each map of relevance and S/N ratio, we find 498 that there is generally a strong positive correlation $(0.71 \le \rho \le 0.77)$ between the LRP 499 indicator patterns and the S/N ratios for temperature, and a moderate positive corre-500 lation $(0.30 \le \rho \le 0.56)$ for precipitation. The exact correlation coefficients between 501 each map are displayed in the subtitles for Figures 6e-h and 7e-h. 502

Given that precipitation only contributes a small amount of relevance compared 503 to temperature, it is perhaps unsurprising that there are several regions where the S/N 504 ratio for precipitation is high, but the relevance is low (e.g. Alaska in JJA, Figure 7c,g 505 or South Africa in SON, Figure 7d,h). Most likely, the forced signal of temperature is 506 clear enough that these regions do not add to the predictive skill of the neural networks. 507 Regions also exist where the S/N ratio for temperature is high despite low relevance (e.g. 508 North Africa in DJF, Figure 6a,e), although these are more rare, as hinted by the strong 509 correlation between the temperature maps of S/N ratio and relevance. In contrast, there 510 are fewer regions with high relevance despite low S/N ratios, but they do occur (e.g. SON 511 temperatures in northern South America, Figure 6d,h). These high-relevance, low-S/N 512 ratio regions confirm that the indicator patterns identified by LRP capture more than 513 the local S/N ratio. Some reasons a region/variable/season may be important in terms 514 of LRP, but not in terms of S/N ratio, are: 1) LRP may be identifying places in our data 515 where a signal exists only in the combination of regions/seasons/variables, which would 516 not be captured by this definition of S/N ratio. 2) Since LRP highlights the patterns the 517 neural networks use to predict the correct decade over all other decades, it may be cap-518 turing abrupt nonlinear changes in the climate that are filtered out by the century-long 519 analysis of S/N ratio In the next section, we discuss further applications of neural network-520 derived indicator patterns and task the network with the much harder problem of iden-521 tifying changes in extreme precipitation over the Americas. 522

523 5 Extreme Precipitation over the Americas

We now task the neural networks to predict the year given combinations of sea-524 sons for a single variable: extreme precipitation over the Americas. We choose to shift 525 our focus for a few reasons. First, we wish to demonstrate that this neural network ap-526 proach can be extended to variables that have considerable noise (like extreme precip-527 itation, see Figure S8), and datasets that do not cover the globe. Second, extreme pre-528 cipitation has major implications for human health (Ali et al., 2019; Eekhout et al., 2018; 529 Rosenzweig et al., 2002) but there is considerable disagreement between climate mod-530 els in its signal (Figure S8). This neural network approach can be used to identify agreed-531 upon patterns despite climate model spread. Further in this section, we will demonstrate 532 that LRP maps can be used to investigate climate model differences and better under-533 stand the time evolution of the forced response. 534

The extreme precipitation signal is not as pronounced as the temperature signal, and using the Americas rather than the full globe limits the amount of unique information in the input field. Nevertheless, the neural networks are still able to detect patterns

-16-

of forced change. Figure 8 depicts the years predicted by one neural network trained on 538 seasonal-mean extreme precipitation. As in Figure 3, the neural network is unable to ac-539 curately predict the year given CMIP6 data prior to the TOE around 2010, whereafter 540 the predicted years generally follow the 1:1 line with the truth years, indicating that the 541 neural network has identified reliable indicators of change for this period. All Pearson 542 correlations of the actual years with the predicted years for extreme precipitation in ob-543 servations are positive (r ≈ 0.4), demonstrating that the indicator patterns found in cli-544 mate models can be successfully applied to observations (Figure 4). These correlations 545 are not as strong as those for mean precipitation observations, due in part to the mag-546 nitude of climate model disagreement in extreme precipitation as well as the observa-547 tional dataset used: ERA5. As shown in Figure S6, the correlations of actual with pre-548 dicted years for ERA5 precipitation observations are far smaller than those for GPCP 549 observations. ERA5 tends to perform poorly in remote regions such as northern North 550 America and northwestern South America (Bell et al., 2021), which may be responsible 551 for these low correlations. The correlation between actual years and neural network-predicted 552 years for extreme precipitation observations are explored in much more detail by Madakumbura 553 et al. (2021). 554

To investigate the indicator patterns used by the neural networks to predict the 555 year when the forced signal first emerges from the background noise, we apply LRP to 556 all climate model samples in the training and testing sets for all 100 neural networks at 557 the TOE (using the TOE calculated for each climate model and neural network individ-558 ually, see Figure S9). LRP points to western South America in DJF and British Columbia 559 in MAM and SON as the most relevant regions when the neural networks first detect the 560 forced response (Figure 9a-d). These LRP maps exhibit a more even distribution in rel-561 evance across each region and season than the end-of-the-21st-century LRP maps of global 562 temperature and precipitation (Figures 6a-d, 7a-d). Predicting the year at the TOE, when 563 the signal has just barely emerged from the background climate, likely requires the neu-564 ral networks to use all of the information available to them. 565

Up to this point, we have only considered the mean LRP maps across climate mod-566 els. Since the neural networks are nonlinear by nature, they can identify multiple pat-567 terns that differ between climate models for a given decade. We apply k-means cluster-568 ing to all 3200 LRP maps at the TOE (32 climate models samples, 100 neural networks) 569 to identify two distinct indicator patterns that are being used by the climate models (Fig-570 ure 9e-l, see the supplementary materials for more details on k-means clustering). Tak-571 ing the difference between the mean LRP maps for clusters one and two reveals that the 572 Amazon in JJA is a highly relevant region in cluster one, while western Canada in DJF 573

is a highly relevant region in cluster two (Figure 9m-p). With the sole exception of MPI-574 ESM1-2-HR, all 100 LRP maps for each individual climate model fall cleanly into one 575 cluster or the other, suggesting that there are two distinct ways in which the forced sig-576 nal emerges in the CMIP6 simulations (Figure 10). Interestingly, when k-means is in-577 structed to identify 32 unique clusters within the LRP maps, each cluster contains all 578 100 relevance maps for each of the 32 climate models. In other words, the pathway used 579 by the neural networks to predict the year is unique to each climate model and distin-580 guishable from all other climate models, regardless of whether the climate model sam-581 ples appear in the training or testing sets (further investigated by Labe & Barnes, 2022). 582

In the same way that indicator patterns can differ between models, indicator pat-583 terns are also able to evolve through time (e.g., Barnes et al., 2020; Labe & Barnes, 2021; 584 Madakumbura et al., 2021). Comparing the LRP maps at the TOE (Figure 11a-d) with 585 those at the end of the 21st century (Figure 11e-h) highlights the regions that become 586 more important for predicting the year over time. The difference plots in Figure 11i-l 587 reveal that the neural network learns to focus on Alaska during MAM, JJA, and SON, 588 Greenland in JJA and SON, and Quebec in MAM and SON as the forced response be-589 comes stronger. These regions are more important for predicting the year at the end of 590 the 21st century than the early 21st century. While further exploration is required, there 591 are several reasons a region may become more relevant over time. For example, it may 592 be that the region does not initially have a clear forced signal, but following some abrupt 593 change (e.g. an ice-free Arctic) the forced signal becomes extremely pronounced. It may 594 also be that the region has a signal that is consistently agreed upon by the majority of 595 CMIP6 climate models, and becomes more relevant compared to other regions as climate 596 model projections in those other regions drift apart. These time-varying patterns sup-597 port the idea that combined indicators are effective for identifying dynamically evolv-598 ing patterns of forced change. 599

600 6 Conclusions

Neural networks are powerful tools for identifying patterns of forced change in the 601 climate system. When tasked with predicting the year given climate model simulations 602 of temperature, precipitation, or extreme precipitation, artificial neural networks can learn 603 these patterns of forced change that allow them to distinguish between maps from dif-604 ferent years. In combined fields, such as multiple variables, seasons, or both, the forced 605 response can be detected earlier than in single fields alone. By visualizing the decision-606 making process of the neural networks with an explainability method we extracted re-607 liable, multivariate patterns of forced change. These neural network-derived combined 608

-18-

indicator patterns are complex and nonlinear and capture more than the local signal to-noise ratio. Explainability methods take a huge step towards disentangling the rela tionships learned by neural networks by pointing out what inputs contributed most to
 the final prediction, but they stop short of explaining why.

Expanding on previous work by Barnes et al. (2020), we used k-means clustering 613 in tandem with layer-wise relevance propagation to study the relationships learned by 614 the neural networks. This approach revealed two distinct ways in which the extreme pre-615 cipitation response emerges in CMIP6 data. While combining neural network explain-616 ability methods with other statistical techniques can enhance explanations of neural net-617 work decisions, there is still a large gap between what the neural network has learned 618 and what we can explain post hoc. Some unanswered questions, such as "why does tem-619 perature in Region A combine with precipitation in Region B to improve the signal of 620 the forced response?" may be better answered with a different architectural approach, 621 such as neural network designs that are inherently interpretable and do not require post-622 hoc approaches like LRP (Rudin, 2019). This is a natural next step for future work. The 623 flexibility and accessibility of this framework provide several other future research di-624 rections. Given that this predict-the-year approach can be applied to observational data, 625 one possible extension of this work could involve exploring the observed features of forced 626 change that are consistent with climate model simulations. There is also space for these 627 methods to be used to determine which definitions of seasons are optimal for detecting 628 forced change. While we used meteorological seasons here, there may be more appropri-629 ate definitions, such as unique definitions of the wet and dry seasons, or the shoulder sea-630 sons, that vary between variables and regions. Furthermore, this framework should be 631 expanded to other variables, regions of focus, and climate change scenarios, to identify 632 the combined indicators that best elucidate the forced signal. For example, extreme pre-633 cipitation and extreme drought may combine to capture the increased volatility in pre-634 cipitation extremes that are expected with climate change (O'Gorman, 2015). Further 635 application of this technique to compound climate extremes, such as heat wave inten-636 sity, drought duration, and flood frequency, may reveal that explainable neural networks 637 are useful for assessing societal impacts and improving climate change preparedness. 638



Figure 1. Schematic of the fully connected neural network architecture. Inputs from multiple maps of data are flattened into an input layer vector (size of the input layer ranges from 948 to 22,752). These inputs are fed through two hidden layers with ten nodes each. The neural network is trained to predict the year that the data came from, outputting the likelihood that the input data came from each decade midpoint between 1905 and 2115. This is then converted to a year via fuzzy classification.



Figure 2. Calculation of TOE. The TOE is defined as the earliest year in which a map, and all subsequent maps, permanently exceed the maximum predicted year from the baseline period (1920-1959). The baseline maximum for each model is indicated by the horizontal lines, the last year that falls below the baseline maximum is circled, and the TOE is indicated by the vertical lines. Sample model 1 (dark red) has a baseline maximum of 1966 and permanently exceeds this threshold in 2028. Sample model 2 (light green) has a baseline maximum of 1981 and permanently exceeds this threshold in 1989. Thus, the TOE for sample models 1 and 2 are estimated as 2028 and 1989, respectively.



Figure 3. Neural network output for temperature and precipitation. Year predicted by the neural network (y-axis) versus the truth year (x-axis) for temperature (a, d, g), precipitation (b, e, h), and temperature and precipitation combined (c, f, i). Input maps include annualmean data (a, b, c), seasonal-mean data (d, e, f), and monthly-mean data (g, h, i). Testing data is shown in color and observations are shown in white.



Figure 4. Correlation of actual years with predicted years for observations. Pearson correlations of the actual years with the years predicted by 100 trained neural networks given observations of temperature, precipitation, and extreme precipitation. Correlations were computed for all years beginning in 1980 where observational data exists for all variables. The box plots indicate the first, second, and third quartile statistics, and the whiskers denote 1.5 times the interquartile range, or the minumum/maximum value, whichever is less extreme. Outliers are excluded for clarity, but can be found in Figures S5 and S6.



Figure 5. Mean TOE for each input field. Comparison of the mean time of emergence identified by neural networks trained on annual-mean (a), seasonal-mean (b), and monthly-mean (c) input fields, and neural networks trained on temperature (d), precipitation (e), and temperature and precipitation combined (f). 100 neural networks with different train-test splits were trained for each input field. Each dot represents the mean TOE for all climate models in the testing set for a single trained neural network, ranked from earliest to latest. Note the change in the y-axes between panels, and that the TOE results for each set of neural networks appear once in the panels a, b, and c, and once in d, e, and f.



Figure 6. Combined indicator patterns of the forced response (temperature). Average temperature LRP results for the seasonal-mean combined neural networks (left, in yellow) and S/N ratio (right, in green) for 2090. Darker shading indicates regions of temperature that are more relevant for the neural network's prediction or have a higher S/N ratio. The Spearman's rank correlation (ρ) between corresponding maps of relevance and S/N ratio are shown in the subtitles of panels e-h.



Figure 7. Combined indicator patterns of the forced response (precipitation). Average precipitation LRP results for the seasonal-mean combined neural networks (left, in yellow) and S/N ratio (right, in green) for 2090. Darker shading indicates regions of precipitation that are more relevant for the neural network's prediction or have a higher S/N ratio. The Spearman's rank correlation (ρ) between corresponding maps of relevance and S/N ratio are shown in the subtitles of panels e-h.



Figure 8. Neural network output for extreme precipitation. Year predicted by the neural network (y-axis) versus the truth year (x-axis) given seasonal-mean maps of extreme precipitation. Testing data is shown in pink and observations are shown in white.



Figure 9. Relevance map clusters at the TOE for extreme precipitation. Average LRP results for: extreme precipitation at the TOE (a-d), each cluster identified by k-means (e-h, i-l), and the difference between the clusters (m-p). In panels a-l, darker shading indicates regions of extreme precipitation that are more relevant for the neural networks' prediction of the year at the TOE. In panels m-p, blue shading indicates the regions that are more relevant in cluster 1, while red shading indicates the regions that are more relevant in cluster 2. Note that panels a-d are identical to panels a-d in Figure 11.



Figure 10. Climate models in each relevance map cluster at the TOE. The number of times each climate model appears in each cluster when k-means is applied to the maps of relevance at the TOE for 100 ANNs trained on extreme precipitation over the Americas. Only the relevance maps for MPI-ESM1-2-HR appear in both clusters. All other relevance maps for each climate model are found in one cluster or the other.



Figure 11. Time evolution of extreme precipitation relevance. Average LRP results at the TOE (a-d), 2090 (e-h), and the difference between (i-l). Darker shading in panels a-h high-lights regions that were more relevant for the neural networks' prediction of the year. In panels i-l, red shading indicates regions where the relevance has increased over time, while blue shading indicates regions where the relevance has decreased over time. Note that panels a-d are identical to panels a-d in Figure 9.

639 Acknowledgments

We thank the two anonymous reviewers, the Editor, Maria Rugenstein, and Jessica Witt 640 for their constructive feedback, which substantially improved our study. I would also like 641 to thank the members of the Barnes Group for all the science discussions that brought 642 about new ideas, and my parents and siblings for their support through the pandemic. 643 This material is based upon work supported by the U.S. Department of Energy, Office 644 of Science, Office of Advanced Scientific Computing Research, Department of Energy Com-645 putational Science Graduate Fellowship under Award Number DE-SC0020347, and by 646 NOAA MAPP grant NA19OAR4310289. We acknowledge the World Climate Research 647 Programme's Working Group on Coupled Modelling, which is responsible for CMIP, and 648 we thank the climate modeling groups for producing and making available their model 649 output. For CMIP the U.S. Department of Energy's Program for Climate Model Diag-650 nosis and Intercomparison provides coordinating support and led development of soft-651 ware infrastructure in partnership with the Global Organization for Earth System Sci-652 ence Portals. 653

⁶⁵⁴ Data Availability Statement

All data used in this study is publicly available and referenced throughout the paper.

- The CMIP6 simulations used in this study can be via the Earth System Grid Federa-
- tion (ESGF, https://esgf-node.llnl.gov/projects/cmip6/). Monthly temperature obser-
- vations are available through Berkeley Earth (http://berkeleyearth.org/data/). Global
- ⁶⁵⁹ Precipitation Climatology Project monthly global precipitation fields are available through
- the NOAA Physical Sciences Laboratory (https://psl.noaa.gov/data/gridded/data.gpcp.html).
- Monthy, daily, and sub-daily precipitation reanalyses were provided by the European Cen-
- tre for Medium-Range Weather Forecasts (ERA5: https://www.ecmwf.int/en/forecasts/datasets/reanalysis-
- datasets/era5) and the National Center for Atmospheric Research (JRA55: https://climatedataguide.ucar.edu/clim
- data/jra-55). Python code used in this work has been made publicly available at https://github.com/jaminrader/l

665 References

670

- Abiodun, O. I., Jantan, A., Omolara, A. E., Dada, K. V., Mohamed, N. A., & Ar-
- shad, H. (2018, November). State-of-the-art in artificial neural network applications: A survey. *Heliyon*, 4(11), e00938. doi: 10.1016/j.heliyon.2018.e00938
- Adler, R. F., Sapiano, M., Huffman, G. J., Wang, J., Gu, G., Bolvin, D., ... Shin,
 - D.-B. (2018, April). The global precipitation climatology project (GPCP)
- monthly analysis (new version 2.3) and a review of 2017 global precipitation.
- Atmosphere, 9(4). doi: 10.3390/atmos9040138

- Ali, H., Modi, P., & Mishra, V. (2019, September). Increased flood risk in indian
 sub-continent under the warming climate. Weather and Climate Extremes, 25,
 100212. doi: 10.1016/j.wace.2019.100212
- Bach, S., Binder, A., Montavon, G., Klauschen, F., Müller, K.-R., & Samek, W.
 (2015, July). On Pixel-Wise explanations for Non-Linear classifier decisions
 by Layer-Wise relevance propagation. *PLoS One*, 10(7), e0130140. doi:
 10.1371/journal.pone.0130140
- Barnes, E. A., Hurrell, J. W., Ebert-Uphoff, I., Anderson, C., & Anderson, D. (2019,
 November). Viewing forced climate patterns through an AI lens. *Geophys. Res. Lett.*, 46(22), 13389–13398. doi: 10.1029/2019gl084944
- Barnes, E. A., Toms, B., Hurrell, J. W., Ebert-Uphoff, I., Anderson, C., & Anderson, D. (2020, September). Indicator patterns of forced change learned by an
 artificial neural network. J. Adv. Model. Earth Syst., 12(9), e2020MS002195.
 doi: 10.1029/2020ms002195
- Barnett, T. P., Pierce, D. W., Hidalgo, H. G., Bonfils, C., Santer, B. D., Das, T.,
 Dettinger, M. D. (2008, February). Human-induced changes in the hydrology of the western united states. *Science*, 319(5866), 1080–1083. doi:
 10.1126/science.1152538
- Bell, B., Hersbach, H., Simmons, A., Berrisford, P., Dahlgren, P., Horányi, A., ...
- Thépaut, J. (2021). The ERA5 global reanalysis: Preliminary extension to
 1950. Quart J Roy Meteor Soc.
- Bindoff, N. L., Stott, P. A., AchutaRao, K. M., Allen, M. R., Gillett, N., Gutzler,
- D., ... Zhang, X. (2013). Chapter 10 detection and attribution of climate
 change: From global to regional. In *Climate change 2013: The physical sci- ence basis. IPCC working group I contribution to AR5.* Cambridge, United
 Kingdom: Cambridge University Press.
- Bonfils, C. J. W., Santer, B. D., Fyfe, J. C., Marvel, K., Phillips, T. J., & Zimmer-
- man, S. R. H. (2020, July). Human influence on joint changes in temperature,
 rainfall and continental aridity. *Nat. Clim. Chang.*, 10(8), 726–731. doi:
 10.1038/s41558-020-0821-1
- Brenowitz, N. D., & Bretherton, C. S. (2018, June). Prognostic validation of a
 neural network unified physics parameterization. *Geophys. Res. Lett.*, 45(12),
 6289–6298. doi: 10.1029/2018gl078510
- Chen, C., Zhang, P., Zhang, H., Dai, J., Yi, Y., Zhang, H., & Zhang, Y. (2020,
 March). Deep learning on Computational-Resource-Limited platforms: A
 survey. Mobile Information Systems, 2020. doi: 10.1155/2020/8454327
- ⁷⁰⁹ Cui, L., Wang, L., Qu, S., Singh, R. P., Lai, Z., & Yao, R. (2019, April). Spa-

710	tiotemporal extremes of temperature and precipitation during 1960–2015 in
711	the yangtze river basin (china) and impacts on vegetation dynamics. Theor.
712	Appl. Climatol., $136(1)$, 675–692. doi: 10.1007/s00704-018-2519-0
713	Dai, A., Lin, X., & Hsu, KL. (2007, October). The frequency, intensity, and di-
714	urnal cycle of precipitation in surface and satellite observations over low- and
715	mid-latitudes. Clim. Dyn., 29 (7-8), 727–744. doi: 10.1007/s00382-007-0260-y
716	Donat, M. G., Alexander, L. V., Herold, N., & Dittus, A. J. (2016, October). Tem-
717	perature and precipitation extremes in century-long gridded observations,
718	reanalyses, and atmospheric model simulations. J. Geophys. Res., 121(19),
719	11,174–11,189. doi: 10.1002/2016jd025480
720	Eekhout, J. P. C., Hunink, J. E., Terink, W., & de Vente, J. (2018, April).
721	Why increased extreme precipitation under climate change negatively af-
722	fects water security. Hydrol. Earth Syst. Sci. Discuss., $22(11)$, 1–16. doi:
723	10.5194/hess-2018-161
724	Eyring, V., Bony, S., Meehl, G. A., Senior, C. A., Stevens, B., Stouffer, R. J., &
725	Taylor, K. E. (2016, May). Overview of the coupled model intercomparison
726	project phase 6 (CMIP6) experimental design and organization. $Geosci. Model$
727	Dev., $9(5)$, 1937–1958. doi: 10.5194/gmd-9-1937-2016
728	Field, C. B., Barros, V. R., Mastrandrea, M. D., Mach, K. J., Abdrabo, M. AK.,
729	Adger, N., Yohe, G. W. (2014). Summary for policymakers. In <i>Climate</i>
730	change 2014: Impacts, adaptation, and vulnerability. part a: Global and sec-
731	toral aspects. contribution of working group II to the fifth assessment report of
732	the intergovernmental panel on climate change (pp. 1–32). Cambridge, United
733	Kingdom and New York, NY, USA: Cambridge University Press.
734	Fischer, E. M., & Knutti, R. (2012, September). Robust projections of combined hu-
735	midity and temperature extremes. Nat. Clim. Chang., $3(2)$, 126–130. doi: 10
736	.1038/nclimate1682
737	Gaetani, M., Janicot, S., Vrac, M., Famien, A. M., & Sultan, B. (2020, May). Ro-
738	bust assessment of the time of emergence of precipitation change in west africa.
739	Sci. Rep., $10(1)$, 7670. doi: 10.1038/s41598-020-63782-2
740	Gettelman, A., Gagne, D. J., Chen, CC., Christensen, M. W., Lebo, Z. J., Mor-
741	rison, H., & Gantos, G. (2021, February). Machine learning the warm
742	rain process. J. Adv. Model. Earth Syst., $13(2)$, e2020MS002268. doi:
743	$10.1029/2020 \mathrm{ms} 002268$
744	Hawkins, E., Frame, D., Harrington, L., Joshi, M., King, A., Rojas, M., & Sutton,
745	R. (2020, March). Observed emergence of the climate change signal: From
746	the familiar to the unknown. <i>Geophys. Res. Lett.</i> , $47(6)$, e2019GL086259. doi:

747	10.1029/2019gl 086259
748	Hawkins, E., & Sutton, R. (2009, August). The potential to narrow uncertainty in
749	regional climate predictions. Bull. Am. Meteorol. Soc., 90(8), 1095–1108. doi:
750	10.1175/2009BAMS2607.1
751	Hawkins, E., & Sutton, R. (2012, January). Time of emergence of climate signals.
752	Geophys. Res. Lett., 39(1). doi: 10.1029/2011gl050087
753	Hersbach, H., Bell, B., Berrisford, P., Hirahara, S., Horányi, A., Muñoz-Sabater, J.,
754	Jean-Noël Thépaut (2020, July). The ERA5 global reanalysis. Quart. J.
755	Roy. Meteor. Soc., $146(730)$, 1999–2049. doi: 10.1002/qj.3803
756	Kim, YH., Min, SK., Zhang, X., Sillmann, J., & Sandstad, M. (2020, September).
757	Evaluation of the CMIP6 multi-model ensemble for climate extreme indices.
758	Weather and Climate Extremes, 29, 100269. doi: 10.1016/j.wace.2020.100269
759	Labe, Z. M., & Barnes, E. A. (2021, June). Detecting climate signals using explain-
760	able AI with single-forcing large ensembles. J. Adv. Model. Earth Syst., $13(6)$,
761	e2021MS002464. doi: $10.1029/2021ms002464$
762	Labe, Z. M., & Barnes, E. A. (2022). Comparison of climate model large ensem-
763	bles with observations in the arctic using simple neural networks. doi: $10.1002/$
764	essoar.10510977.1
765	Lagerquist, R., McGovern, A., & Gagne, D. J., II. (2019, August). Deep learning for
766	spatially explicit prediction of Synoptic-Scale fronts. Weather Forecast., $34(4)$,
767	1137–1160. doi: 10.1175/WAF-D-18-0183.1
768	Lapuschkin, S. (2019). Opening the machine learning black box with layer-wise
769	relevance propagation (Doctoral dissertation, Technischen Universität Berlin,
770	Berlin, Germany). doi: 10.14279/DEPOSITONCE-7942
771	Lee, Y., Kummerow, C. D., & Ebert-Uphoff, I. (2021, April). Applying machine
772	learning methods to detect convection using geostationary operational environ-
773	mental satellite-16 (GOES-16) advanced baseline imager (ABI) data. Atmos.
774	Meas. Tech., 14(4), 2699–2716. doi: 10.5194/amt-14-2699-2021
775	Li, J., Thompson, D. W. J., Barnes, E. A., & Solomon, S. (2017, December). Quan-
776	tifying the lead time required for a linear trend to emerge from natural climate
777	variability. J. Clim., 30(24), 10179–10191. doi: 10.1175/JCLI-D-16-0280.1
778	Lu, E., Chen, H., Tu, J., Song, J., Zou, X., Zhou, B., Jiang, Z. (2015, De-
779	cember). The nonlinear relationship between summer precipitation in
780	china and the sea surface temperature in preceding seasons: A statistical
781	demonstration. J. Geophys. Res. D: Atmos., $120(23)$, $12,027-12,036$. doi:
782	10.1002/2015 JD024030

⁷⁸³ Madakumbura, G. D., Thackeray, C. W., Norris, J., Goldenson, N., & Hall, A.

(2021, July). Anthropogenic influence on extreme precipitation over global 784 land areas seen in multiple observational datasets. Nat. Commun., 12(1), 3944. 785 doi: 10.1038/s41467-021-24262-x 786 Maher, N., Milinski, S., & Ludwig, R. (2021). Large ensemble climate model simula-787 tions: introduction, overview, and future prospects for utilising multiple types 788 of large ensemble. Earth System Dynamics, 12(2), 401-418. 789 Mahony, C. R., & Cannon, A. J. (2018, February). Wetter summers can intensify 790 departures from natural variability in a warming climate. Nat. Commun., 9(1), 791 783. doi: 10.1038/s41467-018-03132-z 792 Mamalakis, A., Ebert-Uphoff, I., & Barnes, E. A. (2021). Neural network attribu-793 tion methods for problems in geoscience: A novel synthetic benchmark dataset. 794 arXiv. Retrieved from http://arxiv.org/abs/2103.10005 795 Mankin, J. S., Lehner, F., Coats, S., & McKinnon, K. A. The (2020, October). 796 value of initial condition large ensembles to robust adaptation decision-making. 797 Earths Future, 8(10). doi: 10.1029/2020ef001610798 Marvel, K., & Bonfils, C. (2013, November). Identifying external influences on global 799 precipitation. Proc. Natl. Acad. Sci. U. S. A., 110(48), 19301-19306. doi: 10 800 .1073/pnas.1314382110 801 Montavon, G., Binder, A., Lapuschkin, S., Samek, W., & Müller, K.-R. (2019).802 Layer-Wise relevance propagation: An overview. In W. Samek, G. Montavon, 803 A. Vedaldi, L. K. Hansen, & K.-R. Müller (Eds.), Explainable AI: Interpreting, 804 explaining and visualizing deep learning (Vol. 11700, pp. 193–209). Cham: 805 Springer International Publishing. doi: 10.1007/978-3-030-28954-6_10 806 Trend analysis of climate time series: A review of Mudelsee, M. (2019, March). 807 methods. Earth-Sci. Rev., 190, 310-322. doi: 10.1016/j.earscirev.2018.12.005 808 Newell, G. J., & Lee, B. (1981, May). Ridge regression: An alternative to multiple 809 linear regression for highly correlated data. J. Food Sci., 46(3), 968-969. doi: 810 10.1111/j.1365-2621.1981.tb15400.x 811 North, G. R., & Stevens, M. J. (1998, April). Detecting climate signals in 812 the surface temperature record. J. Clim., 11(4), 563-577. doi: 10.1175/ 813 1520-0442(1998)011(0563:DCSITS)2.0.CO;2 814 O'Gorman, P. A. (2015). Precipitation extremes under climate change. Curr Clim 815 Change Rep, 1(2), 49-59. doi: 10.1007/s40641-015-0009-3 816 O'Neill, B. C., Tebaldi, C., van Vuuren, D. P., Eyring, V., Friedlingstein, P., Hurtt, 817 G., ... Sanderson, B. M. (2016, September). The scenario model intercompari-818 son project (ScenarioMIP) for CMIP6. Geoscientific Model Development, 9(9), 819 3461-3482. doi: 10.5194/gmd-9-3461-2016 820

- Rohde, R. A., & Hausfather, Z. (2020, December). The berkeley earth land/ocean
 temperature record. *Earth Syst. Sci. Data*, 12(4), 3469–3479. doi: 10.5194/
 essd-12-3469-2020
- Rosenzweig, C., Tubiello, F. N., Goldberg, R., Mills, E., & Bloomfield, J. (2002, October). Increased crop damage in the US from excess precipitation under climate change. *Glob. Environ. Change*, 12(3), 197–202. doi: 10.1016/S0959
 -3780(02)00008-0
- Rudin, C. (2019, May). Stop explaining black box machine learning models for high
 stakes decisions and use interpretable models instead. Nature Machine Intelli gence, 1(5), 206-215. doi: 10.1038/s42256-019-0048-x
- Sanderson, B. M., Oleson, K. W., Strand, W. G., Lehner, F., & O'Neill, B. C. (2018,
 February). A new ensemble of GCM simulations to assess avoided impacts in
 a climate mitigation scenario. *Clim. Change*, 146(3), 303–318. doi: 10.1007/
 s10584-015-1567-z
- Santer, B. D., Fyfe, J. C., Solomon, S., Painter, J. F., Bonfils, C., Pallotta, G., &
 Zelinka, M. D. (2019, October). Quantifying stochastic uncertainty in detection time of human-caused climate signals. *Proc. Natl. Acad. Sci. U. S. A.*,
 116(40), 19821–19827. doi: 10.1073/pnas.1904586116
- Santer, B. D., Mears, C., Doutriaux, C., Caldwell, P., Gleckler, P. J., Wigley,
- T. M. L., ... Wentz, F. J. (2011, November). Separating signal and noise
 in atmospheric temperature changes: The importance of timescale. J. Geophys. *Res.*, 116(D22). doi: 10.1029/2011jd016263
- Santer, B. D., Taylor, K. E., Wigley, T. M. L., Johns, T. C., Jones, P. D., Karoly,
- D. J., ... Tett, S. (1996, July). A search for human influences on the
 thermal structure of the atmosphere. *Nature*, 382(6586), 39–46. doi:
 10.1038/382039a0
- Scaife, A. A., & Smith, D. (2018, July). A signal-to-noise paradox in climate science.
 npj Climate and Atmospheric Science, 1(1), 1–8. doi: 10.1038/s41612-018-0038
 -4
- Schneider, T., & Held, I. M. (2001, February). Discriminants of Twentieth-Century
 changes in earth surface temperatures. J. Clim., 14(3), 249–254. doi: 10.1175/
 1520-0442(2001)014(0249:LDOTCC)2.0.CO;2
- Schulzweida, U. (2019). CDO user guide (version 1.9. 6). Max Planck Institute for
 Meteorology: Hamburg, Germany.
- Silva, S. J., Ma, P.-L., Hardin, J. C., & Rothenberg, D. (2021, May). Physically regularized machine learning emulators of aerosol activation. *Geosci. Model Dev.*,
 14(5), 3067–3077. doi: 10.5194/gmd-14-3067-2021

858	Sippel, S., Meinshausen, N., Fischer, E. M., Székely, E., & Knutti, R. (2020, Jan-
859	uary). Climate change now detectable from any single day of weather at global
860	scale. Nat. Clim. Chang., $10(1)$, 35–41. doi: 10.1038/s41558-019-0666-7
861	Solomon, A., & Newman, M. (2012, July). Reconciling disparate twentieth-century
862	Indo-Pacific ocean temperature trends in the instrumental record. Nat. Clim.
863	Chang., $2(9)$, 691–699. doi: 10.1038/nclimate1591
864	Solow, A. R. (1987, October). Testing for climate change: An application of the
865	Two-Phase regression model. J. Appl. Meteorol. Climatol., 26(10), 1401–1405.
866	doi: $10.1175/1520-0450(1987)026\langle 1401: TFCCAA \rangle 2.0.CO; 2$
867	Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., & Salakhutdinov, R.
868	(2014). Dropout: a simple way to prevent neural networks from overfitting.
869	J. Mach. Learn. Res., 15(1), 1929–1958.
870	Swain, D. L., Wing, O. E. J., Bates, P. D., Done, J. M., Johnson, K. A., &
871	Cameron, D. R. (2020, November). Increased flood exposure due to climate
872	change and population growth in the united states. Earths Future, $\mathcal{S}(11)$. doi:
873	10.1029/2020 ef 001778
874	Tabari, H., & Willems, P. (2018, March). Seasonally varying footprint of climate
875	change on precipitation in the middle east. Sci. Rep., $\delta(1)$, 4435. doi: 10.1038/
876	s41598-018-22795-8
877	Toms, B. A., Barnes, E. A., & Ebert-Uphoff, I. (2020, September). Physically
878	interpretable neural networks for the geosciences: Applications to earth sys-
879	tem variability. J. Adv. Model. Earth Syst., $12(9)$, e2019MS002002. doi:
880	10.1029/2019 ms 002002
881	USGCRP. (2018). Impacts, risks, and adaptation in the united states: Fourth na-
882	tional climate assessment, volume II (Tech. Rep.). Washington, DC, USA:
883	U.S. Global Change Research Program. doi: 10.7930/NCA4.2018
884	Weyn, J. A., Durran, D. R., & Caruana, R. (2020, September). Improving data-
885	driven global weather prediction using deep convolutional neural networks
886	on a cubed sphere. J. Adv. Model. Earth Syst., $12(9)$, e2020MS002109. doi:
887	10.1029/2020ms002109
888	
	Wills, R. C., Battisti, D. S., Armour, K. C., Schneider, T., & Deser, C. (2020, Oc-
889	Wills, R. C., Battisti, D. S., Armour, K. C., Schneider, T., & Deser, C. (2020, October). Pattern recognition methods to separate forced responses from inter-
889 890	 Wills, R. C., Battisti, D. S., Armour, K. C., Schneider, T., & Deser, C. (2020, October). Pattern recognition methods to separate forced responses from internal variability in climate model ensembles and observations. J. Clim., 33(20),
889 890 891	 Wills, R. C., Battisti, D. S., Armour, K. C., Schneider, T., & Deser, C. (2020, October). Pattern recognition methods to separate forced responses from internal variability in climate model ensembles and observations. J. Clim., 33(20), 8693–8719. doi: 10.1175/JCLI-D-19-0855.1
889 890 891 892	 Wills, R. C., Battisti, D. S., Armour, K. C., Schneider, T., & Deser, C. (2020, October). Pattern recognition methods to separate forced responses from internal variability in climate model ensembles and observations. J. Clim., 33(20), 8693–8719. doi: 10.1175/JCLI-D-19-0855.1 Wills, R. C., Schneider, T., Wallace, J. M., Battisti, D. S., & Hartmann, D. L.
889 890 891 892 893	 Wills, R. C., Battisti, D. S., Armour, K. C., Schneider, T., & Deser, C. (2020, October). Pattern recognition methods to separate forced responses from internal variability in climate model ensembles and observations. J. Clim., 33(20), 8693–8719. doi: 10.1175/JCLI-D-19-0855.1 Wills, R. C., Schneider, T., Wallace, J. M., Battisti, D. S., & Hartmann, D. L. (2018, March). Disentangling global warming, multidecadal variability, and

Non-

doi:

10.1002/2017 GL076327895 Wu, J., Chen, X., Yao, H., Gao, L., Chen, Y., & Liu, M. (2017, August). 896 linear relationship of hydrological drought responding to meteorological 897 drought and impact of a large reservoir. J. Hydrol., 551, 495–507. 898 10.1016/j.jhydrol.2017.06.029 899 Zadeh, L. A. (1965). Information and control. Fuzzy Sets and Systems, 8(3), 338-900 353. doi: 10.1002/joc.1027 901 Zappa, G., Hoskins, B. J., & Shepherd, T. G. (2015, August). Improving climate 902 change detection through optimal seasonal averaging: The case of the north 903

atlantic jet and european precipitation. J. Clim., 28(16), 6381-6397. doi: 904 10.1175/JCLI-D-14-00823.1 905

Supporting Information for "Detection of forced change within combined climate fields using explainable neural networks"

Jamin K. Rader¹, Elizabeth A. Barnes¹, Imme Ebert-Uphoff^{2,3}, Chuck

$Anderson^4$

¹Department of Atmospheric Science, Colorado State University, Fort Collins, CO, USA

²Cooperative Institute for Research in the Atmosphere, Colorado State University, Fort Collins, CO, USA

³Department of Electrical and Computer Engineering, Colorado State University, Fort Collins, CO, USA

 $^4\mathrm{Department}$ of Computer Science, Colorado State University, Fort Collins, CO, USA

Contents of this file

- 1. Text S1: Neural Network Specifications
- 2. Text S2: Selection of Neural Network Hyperparameters
- 3. Text S3: K-means Clustering
- 4. Text S4: Additional Observational Datasets
- 5. Figures S1 to S12
- 6. References

Corresponding author: Jamin K. Rader, (jamin.rader@colostate.edu)

S1 Neural Network Specifications

All units of the neural networks use a rectified linear unit (ReLU) activation function, except for the output layer which uses a soft-max layer to rescale the final outputs of the neural network such that they sum to one. We train the neural networks using the binary cross-entropy loss between the predicted class likelihoods and the correct class membership weights, such that the loss function is minimized when the two are equal. More information on the ReLU activation function, the soft-max layer, and the loss function can be found in sections A1, A2, and A3 of Barnes et al. (2020), respectively.

The neural networks were trained using the Keras Adam optimizer, an adaptive stochastic gradient descent algorithm (Kingma & Ba, 2014). We used a learning rate that started at 0.001 and decayed linearly to 0.0005 over the span of 150 epochs. Although the Adam optimizer is designed to alter the learning rate based on the momentum of training, the decaying learning rate allowed the neural networks to train more quickly with improved performance. Weights and biases were initialized using random values from a normal distribution.

As discussed in Section 3.1, our neural networks are fully connected with two hidden layers and 10 nodes each in each layer. We found that this architecture allowed the neural networks to capture forced change better than a linear model or even simpler architectures, such as neural networks with only one hidden layer or five nodes in each layer (Figure S2). The additional performance offered by more complicated architectures was small and increased the computational resources needed for training. We elected to stick with the simplest model that performed well with minimal computational expense. These neural networks can be trained on standard laptop

or desktop computers in two to ten minutes depending on the input field, making them extremely accessible to those in the climate science community.

As discussed in Section 3.2 and Figure S4, we applied a ridge penalty (L2 regularization) to the input layer (see Barnes et al., 2020). The ridge penalty was selected such that the time of emergence detected by the neural networks was the earliest. All input vectors used a ridge penalty of 0.1, except for seasonal-mean temperature and precipitation combined input vector, for which the TOE was earlier for a ridge penalty of 0.01 (see Figure S4).

Summary of Neural Network Specifications

Number of Hidden Layers	2
Number of Nodes in Each Hidden Layer	10
Hidden Layer Activation Function	ReLU (Rectified Linear Unit)
Output Layer Activation Function	Softmax
Ridge Penalty (applied to the weights of the	0.01 for seasonal-mean temperature and precip
first hidden layer)	0.1 for all other input fields
Loss Function	Binary Cross-entropy
Optimizer	Adam, tensorflow.keras.optimizers.Adam
Learning Rate	Started == s at 0.001, decaying linearly to 0.0005
Number of Epochs	150

S2 Selection of Neural Network Hyperparameters

We explored a range of values for several neural network hyperparameters such as the learning rate (from 10^{-4} to 10^{-1}), the number of epochs (up to 1000), the ridge penalty (from 0 to 1, see Figure S4), and the neural network architecture, where we examined the performance of neural networks with 1, 2, or 3 hidden layers, and 5, 10, 20 or 50 nodes in each hidden layer (see Figure S2). To choose these hyperparameters we employed a strategy similar to leavep-out cross-validation which is commonly used in the atmospheric sciences (Celisse & Robin, 2008). Specifically, we used 10 different train/test splits to explore the hyperparameter space

and optimize the performance of our neural networks. Using 10 different train/test splits, rather than just one, ensures that our hyperparameter selections are not overfitting to any one specific way the climate models can be split into training and testing sets. Once the best hyperparameter choices were made, we then used another 100 train/test splits for the results of this study, all of which differed from the train/test splits used for tuning.

S3 K-means Clustering

Before applying k-means clustering, all LRP maps are converted into binary maps. Every grid point on each LRP map is assigned a one or a zero depending on whether its relevance value is greater than or less than the mean relevance across all maps and grid points. In this way, ones indicate regions of high relevance, and zeros indicate regions of low relevance. K-means clustering is then applied to these binary LRP maps (3200 in total, samples from 32 climate models for 100 neural networks). We used Sci-Kit Learn's sklearn.cluster.KMeans function (version 0.22.1) in Python with 100 different initializations and all other choices were left as default (Pedregosa et al., 2011). The results for K = 2 are shown in the main paper. Using K = 2 identified two clusters that were near-equal in size, and several runs of k-means with different random initial conditions yielded near-identical results. Clustering for K = 3, 4, 5, 6, 7, 8, and 32 was also explored, however the results for three or more clusters were less physically consistent.

S4 Additional Observational Datasets

In addition to the observational datasets in Section 2.2, we also test two additional precipitation observations in Figures S5 and S6. First, we use the European Center for Medium-Range Weather Forecasts' ERA5 global reanalysis (Hersbach et al., 2020) at 6-hour resolution to con-

struct observational monthly mean precipitation fields from 1980 to the present. Second, we use the Japan Meteorological Agency's Japanese 55-year Reanalysis (JRA55; Kobayashi et al., 2015) mean 3-hour precipitation forecasts to construct observational monthly mean precipitation fields from 1959 to the present.



Figure S1. TOE detected by the neural networks given different definitions of season. As in Figure 5d-f, but for each possible three-month combination of seasons. All three definitions lead to similar TOE when neural networks are trained on global maps of temperature or precipitation. When temperature and precipitation are combined, meteorological seasons lead to the earliest detection of forced change.



Figure S2. Skill across various neural network architectures. The mean testing binary cross-entropy loss for ten trained neural networks with different train/test splits for 12 different neural network architectures and one linear model for annual-mean temperature and annual-mean precipitation. The white box indicates the neural network architecture that was used in the main text (two hidden layers, 10 nodes each).



Figure S3. Climate models used for each input variable. Temperature, precipitation, and temperature and precipitation combined used the same 37 CMIP6 climate models. Extreme precipitation fields came from 32 climate models for which daily precipitation fields were available.



Figure S4. TOE and RMSE for various ridge penalties. The sensitivity of RMSE and TOE to the ridge (L2) penalty used for 10 neural networks trained on seasonal-mean maps of (a) temperature, (b) precipitation, and (c) temperature and precipitation combined. Each plot shows the RMSE and TOE for neural networks trained with a ridge penalty of 1, 0.1, and 0.01 (denoted by red circles, blue stars, and orange triangles, respectively). The mean RMSE and TOE for all 10 neural networks are indicated by the horizontal and vertical lines. Each neural network for a given variable/ridge penalty differs only in which climate models were part of the training and testing sets. While a ridge penalty of 0.01 leads to the smallest mean RMSE in all cases, using a higher ridge penalty of 0.1 leads to earlier detection of forced change for temperature and precipitation input vectors. As a result, we choose to use the ridge penalties corresponding to an earlier TOE.



Figure S5. Sensitivity of observational correlations to the source of precipitation observations: temperature and precipitation combined. Pearson correlations of the actual years with the years predicted by 100 trained neural networks given observations of temperature and precipitation. Correlations were computed for all years beginning in 1980 where observational data exists for all variables. The box plots indicate the first, second, and third quartile statistics, and the whiskers denote 1.5 times the interquartile range, or the minumum/maximum value, whichever is less extreme. The observational correlations for seasonal-mean combined neural networks are sensitive to the dataset of choice, as observational correlations are higher for GPCP than ERA5 or JRA55. This is not the case for the annual-mean and monthly-mean combined neural networks, which have approximately the same correlations regardless of the source of the observations. This is because the seasonal-mean combined neural networks rely on precipitation to predict the year, while the annual-mean and monthly-mean combined neural networks do not, as shown in Figure 5.



Figure S6. Sensitivity of observational correlations to the source of precipitation observations: precipitation only. Pearson correlations of the actual years with the years predicted by 100 trained neural networks given observations of precipitation. Correlations were computed for all years beginning in 1980 where observational data exists for all variables. The box plots indicate the first, second, and third quartile statistics, and the whiskers denote 1.5 times the interquartile range, or the minumum/maximum value, whichever is less extreme. The observational correlations are sensitive to the source of precipitation data. Correlations are highest for GPCP, followed by ERA5 and JRA55. The observational correlations for ERA5 seasonal-mean extreme precipitation are similar to those for ERA5 seasonal-mean precipitation.



Figure S7. Time of emergence for seasonal-mean fields. TOE was calculated for each climate model in the testing sets of 100 trained neural networks. Each dot represents five (rounded up) occurrences of the associated TOE year (i.e. one dot represents 1-5 occurrences, two dots represent 6-10 occurrences, and so on). For added clarity, box plots indicate the first, second, and third quartiles of the TOEs for each model, and whiskers denote 1.5 times the interquartile range, or the minimum/maximum point, whichever is less extreme.



Figure S8. Signal and noise for extreme precipitation over the Americas. Plots of S/N ratio, and its components (signal, climate model variability, and internal variability) for extreme precipitation in each season over North and South America. The signal is most clear over the northern-most latitudes. The S/N ratio is below 1.5 in all seasons indicating that there is considerable noise relative to the signal of change.



Extreme Precipitation

Figure S9. Time of emergence for extreme precipitation over the Americas. TOE was calculated for each climate model in the testing sets of 100 trained neural networks. Each dot represents five (rounded up) occurrences of the associated TOE year (i.e. one dot represents 1-5 occurrences, two dots represent 6-10 occurrences, and so on). For added clarity, box plots indicate the first, second, and third quartiles of the TOEs for each model, and whiskers denote 1.5 times the interquartile range, or the minimum/maximum point, whichever is less extreme.



Figure S10. Neural network output for temperature and precipitation (with training data included). Same as Figure 3 with training data included. Year predicted by the neural network (y-axis) versus the truth year (x-axis) for temperature (a, d, g), precipitation (b, e, h), and temperature and precipitation combined (c, f, i). Input maps include annual-mean data (a, b, c), seasonal-mean data (d, e, f), and monthly-mean data (g, h, i). Training data is shown in gray, testing data is shown in color, and observations are shown in white.



Extreme precipitation

Figure S11. Neural network output for extreme precipitation (with training data included). Same as Figure 8 with training data included. Year predicted by the neural network (y-axis) versus the truth year (x-axis) given seasonal-mean maps of extreme precipitation. Training data is shown in gray, testing data is shown in pink, and observations are shown in white.



Figure S12. Learning curves for temperature and precipitation. Binary cross-entropy loss versus epoch of training for the training and testing data for the nine trained neural networks shown in Figure 3, Figure S10.

References

- Barnes, E. A., Toms, B., Hurrell, J. W., Ebert-Uphoff, I., Anderson, C., & Anderson, D. (2020, September). Indicator patterns of forced change learned by an artificial neural network. J. Adv. Model. Earth Syst., 12(9), e2020MS002195. doi: 10.1029/2020ms002195
- Celisse, A., & Robin, S. (2008, January). Nonparametric density estimation by exact leave-p-out cross-validation. *Comput. Stat. Data Anal.*, 52(5), 2350–2368. doi: 10.1016/j.csda.2007.10 .002
- Hersbach, H., Bell, B., Berrisford, P., Hirahara, S., Horányi, A., Muñoz-Sabater, J., ... Jean-Noël Thépaut (2020, July). The ERA5 global reanalysis. Quart. J. Roy. Meteor. Soc., 146(730), 1999–2049. doi: 10.1002/qj.3803
- Kingma, D. P., & Ba, J. (2014, December). Adam: A method for stochastic optimization.
- Kobayashi, S., Ota, Y., Harada, Y., Ebita, A., Moriya, M., Onoda, H., ... Takahashi, K. (2015).
 The JRA-55 reanalysis: General specifications and basic characteristics. 2, 93(1), 5–48.
 doi: 10.2151/jmsj.2015-001
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., ... Duchesnay,
 E. (2011). Scikit-learn: Machine learning in Python. Journal of Machine Learning Research, 12, 2825–2830.