

Explainable AI for Climate Science

Opening the Black Box to Reveal Planet Earth

Dr. Elizabeth A. Barnes
Professor, Dept. of Atmospheric Science
Colorado State University

This research has been funded, in part, by grants
from the NSF, NOAA, DOE, DARPA and with special
thanks to David Wallerstein

SciPy, July 11 2024

A challenging future



Earth system is threatened

Global environmental change has pushed society well outside of humanity's evolutionary experience, and critical ecosystem services are being threatened.

Humans on the move

Unprecedented numbers of people have and will be impacted and seeking to escape from worsening climate hazards.

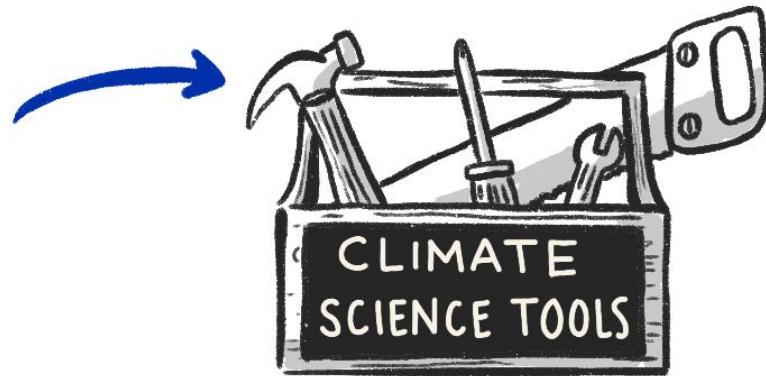


Anticipating what's to come

Predicting human-Earth system futures and change is critical for preparing society as we navigate our turbulent future

AI for climate

MACHINE
LEARNING



AI for climate is not new

But the field's interest and research has exploded in recent years!

1370

MONTHLY WEATHER REVIEW

VOLUME 83

Nonlinear Principal Component Analysis by Neural Networks: Theory and Application to the Lorenz System

ADAM H. MONAHAN

Oceanography Unit, Department of Earth and Ocean Sciences, and
Crisis Points Group, Peter Wall Institute for Advanced Studies, University of British Columbia,
Vancouver, British Columbia, Canada

1998

(Manuscript received 28 October 1998, in final form 7 May 1999)

ABSTRACT

A nonlinear generalization of principal component analysis (PCA), denoted nonlinear principal component analysis (NLPCA), is implemented in a variational framework using a five-layer autoassociative feed-forward neural network. The method is tested on a dataset sampled from the Lorenz attractor, and it is shown that the NLPCA approximations to the attractor in one and two dimensions, explaining 76% and 99.5% of the variance,

1998

(Manuscript received 28 October 1998, in final form 7 May 1999)

ABSTRACT

Two adaptive logic models and a training algorithm for each are described. These models are tested on a meteorological prediction problem with the use of large developmental and test data samples. Discriminant analysis is used for comparison. It is found that ceiling heights at Washington National Airport could be forecast better with the discriminant analysis technique than with the adaptive logic models.

1964

(Manuscript received 12 May 1964, in revised form 25 July 1964)

ABSTRACT

Two adaptive logic models and a training algorithm for each are described. These models are tested on a meteorological prediction problem with the use of large developmental and test data samples. Discriminant analysis is used for comparison. It is found that ceiling heights at Washington National Airport could be forecast better with the discriminant analysis technique than with the adaptive logic models.

1. Introduction

In recent years there has been an increase in the research effort concerning adaptive logic systems. These systems, including their associated hardware, are also known as learning machines or trainable networks. The desire is to create a system that can be trained to give the correct response when it is presented certain input data or stimuli in somewhat the same way a living organism might adjust to its environment or learn from past experience. It is also hoped that the system would be capable of certain generalizations so that when it is presented a data pattern that it has never encountered before it will still tend to give the correct response.

Adaptive logic models can be used for pattern recognition and, since the forecasting of certain weather elements involves, to some extent, the recognition of patterns, an application of adaptive systems to objective weather forecasting is thereby suggested.

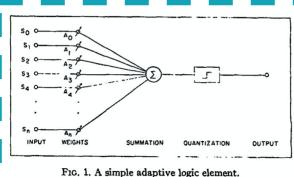
2. Adaptive logic models

One adaptive model that has been described with minor variations by several authors, including Matson,¹ Sebestyen (1962), Widrow (1962), Hu², Duda and Machanki³ and Hu and Root (1964) is described below.

In Fig. 1 each "S" unit represents one of the $n+1$ numerical inputs to the system. Each "A" unit is an

adjustable weight which is altered during the training procedure so that the process of summation, performed in the "S" unit, of $A_i S_i, i=0, 1, 2, \dots, n$, and quantization into a binary output will provide the correct dichotomous classification. S_0 is held constant at +1 while $S_i, i=1, 2, \dots, n$ may be different for different training sample points. This complete element has been called by Widrow (1962) an ADALINE for "adaptive linear neuron."

A training program for this model is indicated in Fig. 2. The value D is a positive constant which permits $\sum_{i=0}^n A_i S_i$ to have a value of D , rather than zero, immediately after an alteration of the A_i for the sample point which occasioned the alteration. After the output is made to agree with the desired output for a particular sample point, the A_i may again be changed for another sample point. This latter alteration may in turn cause the output to be incorrect for the former sample point. Repeated iterations over the data sample can be made and eventually the correct response may be given for all sample points if such distinction is possible with this model.



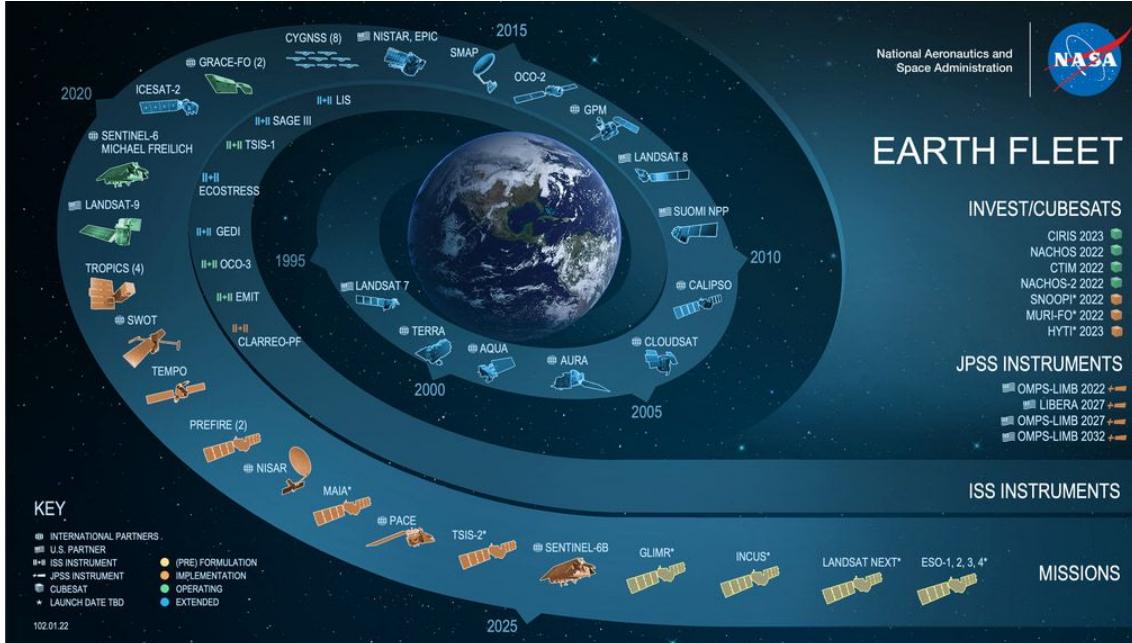
¹ Matson, R. L., 1959: The design and analysis of an adaptive system for statistical classification. S.M. Thesis, Massachusetts Institute of Technology, Cambridge, Mass., 60 pp.

² Hu, C. Y., 1962: A self-adapting binary forecasting system. Tech. Rep. No. 6759-1, Contract No. AF33(616)-7726, Systems Theory Laboratory, Stanford University, Stanford, Calif., 19 pp.

³ Duda, R. O., and P. M. Hart, 1964: An adaptive prediction technique and its application to weather forecasting. Paper presented at the Western Electronic Show and Convention, San Francisco, Calif.

So much data.

Our field has so much (imperfect) data:
ground observations, space-borne
observations, reanalyses, climate model
output, etc...



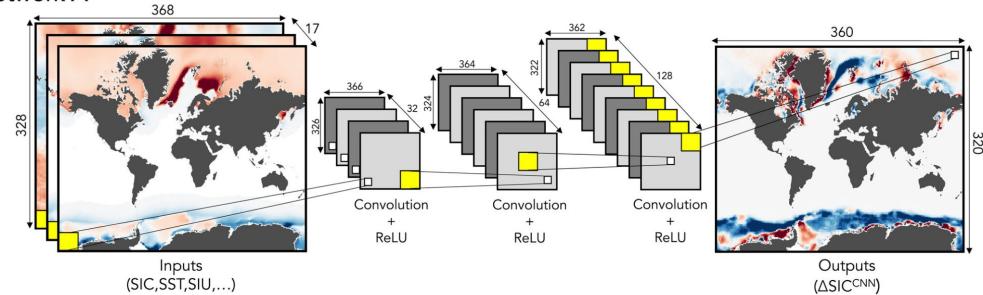
1

ML for post-processing data
[e.g. data compression, data analysis]

Predicting the Errors of Forecast Systems

e.g. Chapman et al. (2019), Cahill et al. (in review), Pan et al. (2021), Gregory et al. (2023)

Network A



Machine Learning: a tool with many uses.

1

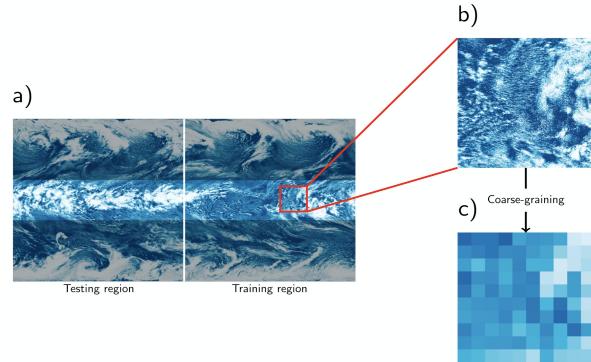
ML for post-processing data
[e.g. data compression, data analysis]

2

ML to improve climate models
[e.g. parameterizations]

Improved Model Parameterizations

e.g. Rasp et al. (2018; PNAS); Schneider et al. (2017; GRL); O'Gorman and Dwyer (2018); Beucler et al. (2020; PRL); Dagon et al (2020); Brenowitz and Bretherton (2018)



Machine Learning: a tool with many uses.

1

ML for post-processing data
[e.g. data compression, data analysis]

2

ML to improve climate models
[e.g. parameterizations]

3

Combine disparate datastreams to
explore complex systems



Machine Learning: a tool with many uses.

1

ML for post-processing data
[e.g. data compression, data analysis]

2

ML to improve climate models
[e.g. parameterizations]

3

Combine disparate datastreams to
explore complex systems

4

Merging observations and model data;
[e.g. transfer learning]

train on climate
model data

fine-tune with limited
observations

predict observations

Machine Learning: a tool with many uses.

1

ML for post-processing data
[e.g. data compression, data analysis]

2

ML to improve climate models
[e.g. parameterizations]

3

Combine disparate datastreams to explore complex systems

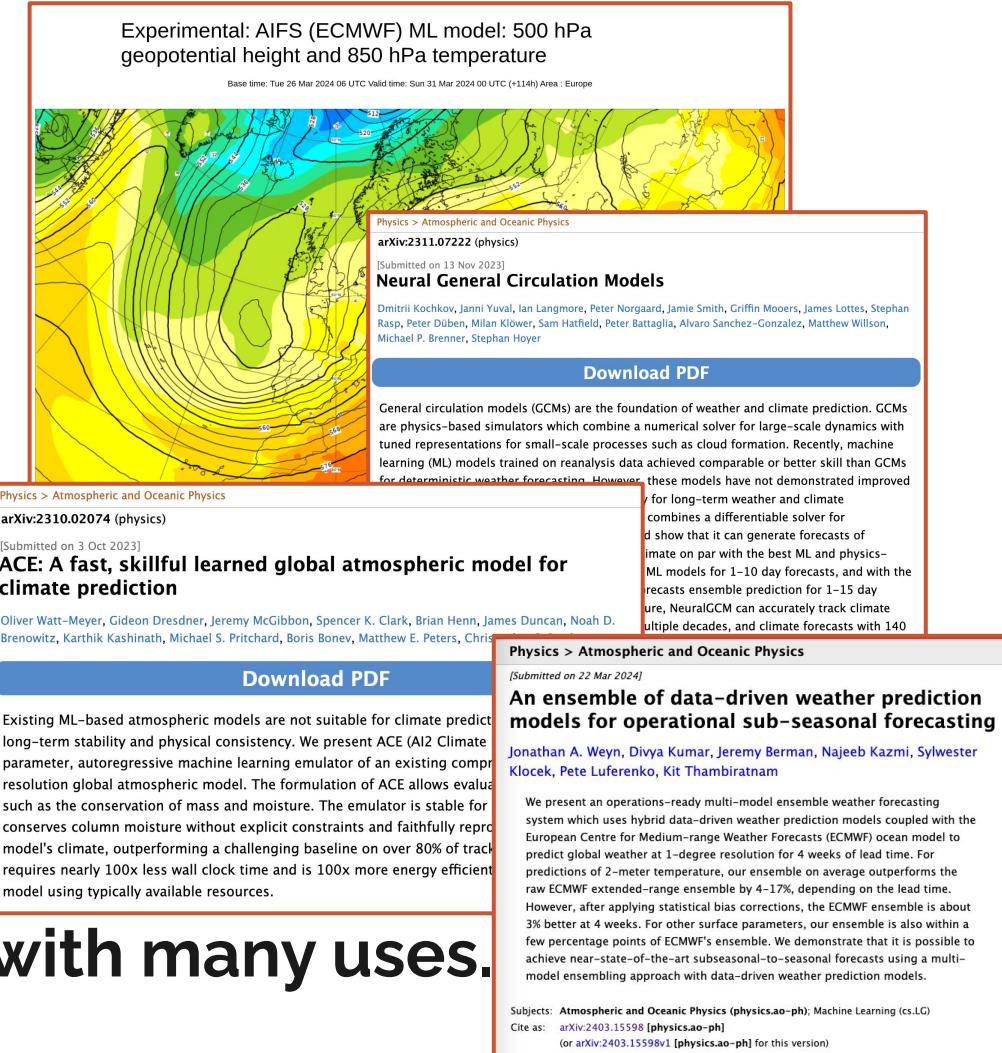
4

Merging observations and model data.
[e.g. transfer learning]

5

Deep-learning weather + climate emulators

Machine Learning: a tool with many uses.



1

ML for post-processing data
[e.g. data compression, data analysis]

2

ML to improve climate models
[e.g. parameterizations]

3

Combine disparate datastreams to
explore complex systems

4

Merging observations and model data.
[e.g. transfer learning]

5

Deep-learning weather + climate
emulators

6

Climate change communication



Machine Learning: a tool with many uses.

ILLUSTRATED BY



1

ML for post-processing data
[e.g. data compression, data analysis]

2

ML to improve climate models
[e.g. parameterizations]

3

Combine disparate datastreams to
explore complex systems

4

Merging observations and model data.
[e.g. transfer learning]

5

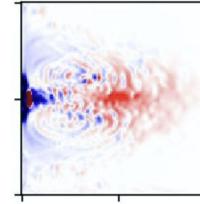
Deep-learning weather + climate
emulators

6

Climate change communication

7

Learn new things!



$$\hat{\mathbf{S}}_{\mathbf{u}}^{BT} \approx \kappa_{BT} \nabla \cdot \begin{pmatrix} \zeta^2 - \zeta D & \zeta \tilde{D} \\ \zeta \tilde{D} & \zeta^2 + \zeta D \end{pmatrix}$$

Equation Discovery

e.g. Zanna & Bolton (2020)

Machine Learning: a tool with many uses.

OUR GOAL:

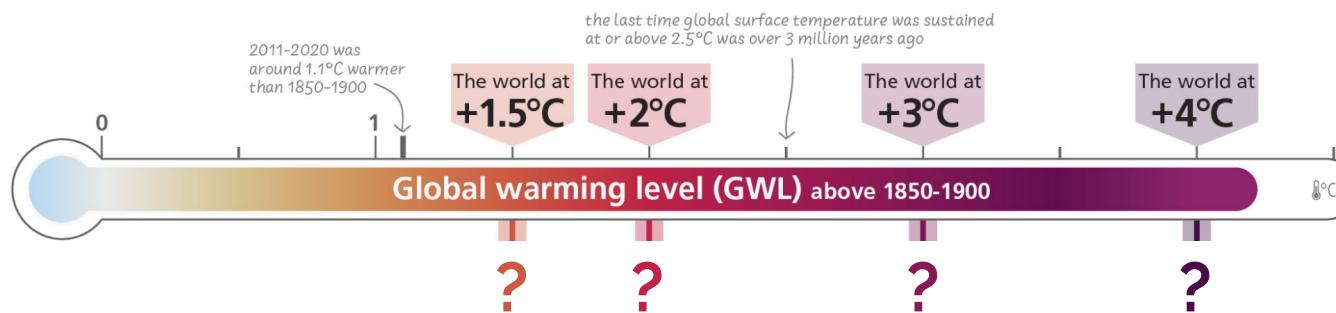
To develop and implement AI tools to leverage imperfect climate models in support of earth system prediction across time and space.

Climate models provide inaccurate, but invaluable “parallel universes” to mine for information

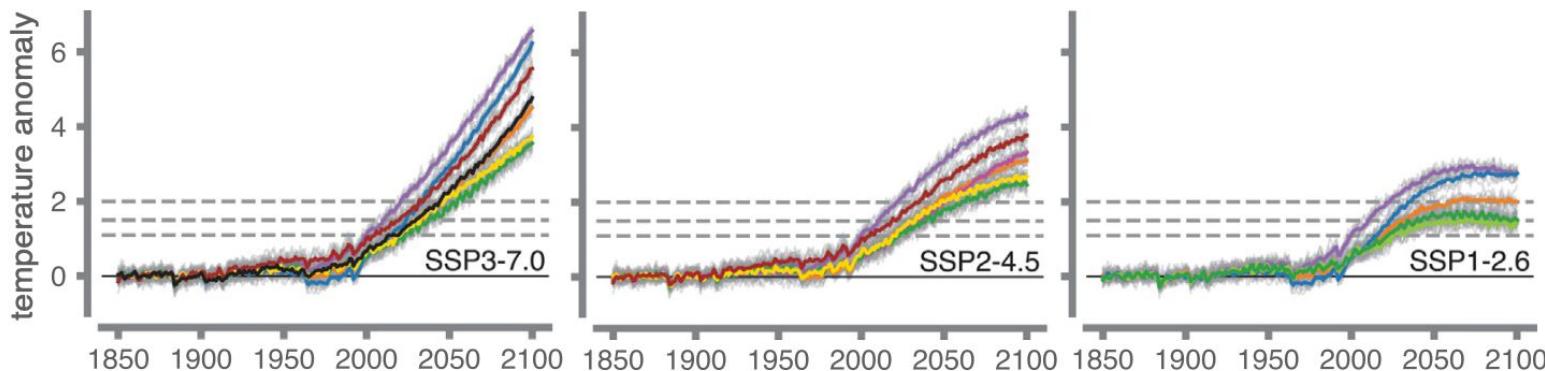


https://futurism.com/parallel-universes-many-worlds-theory

Time Remaining Until Critical Warming Thresholds are Reached

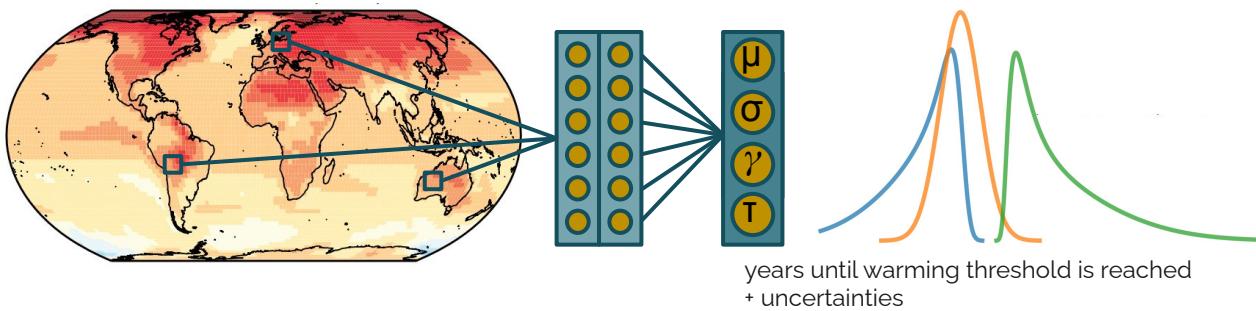


Time Remaining Until Critical Warming Thresholds are Reached



Diffenbaugh & Barnes (2023)

Trained on annual maps from 10 realizations from across multiple climate models

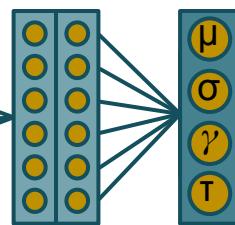
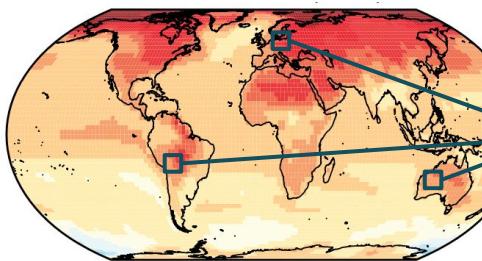


Train neural network to ingest a single annual temperature map and predict the number of years until a warming threshold is reached

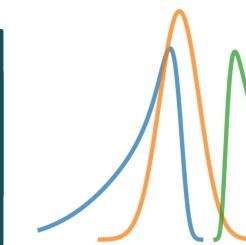


Diffenbaugh & Barnes (2023)

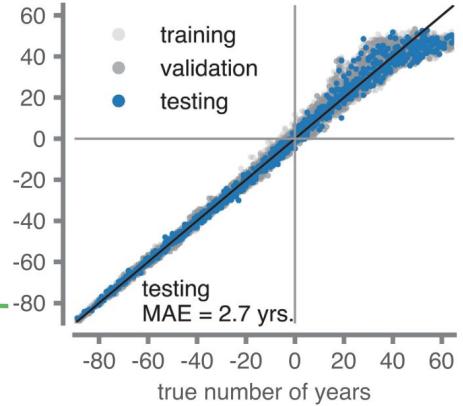
Trained on annual maps from 10 realizations from across multiple climate models



years until warming threshold is reached
+ uncertainties



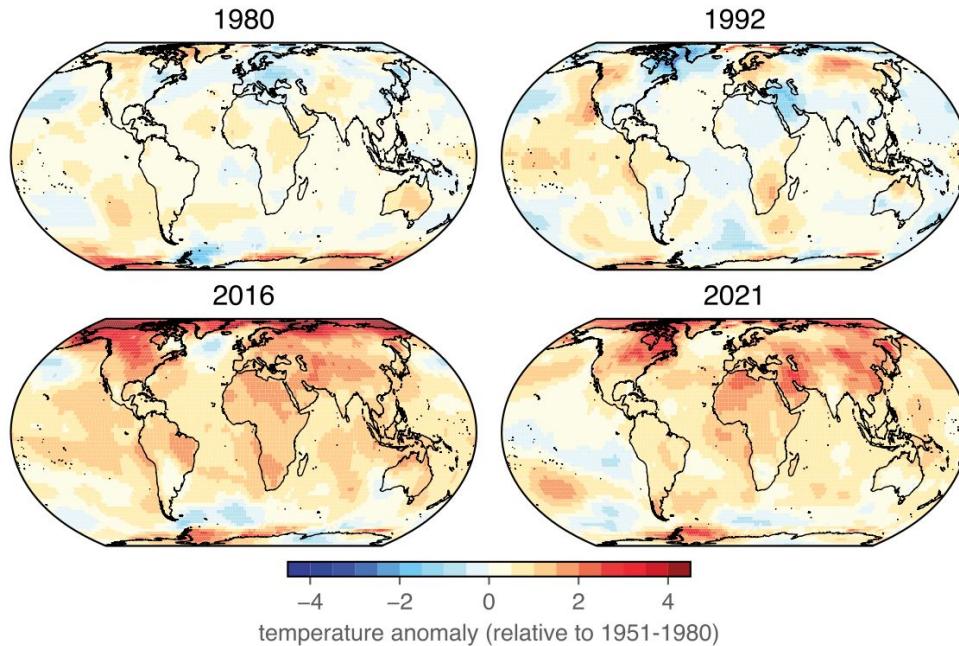
Climate Model Results



Train neural network to ingest a single annual temperature map and predict the number of years until a warming threshold is reached



Diffenbaugh & Barnes (2023)



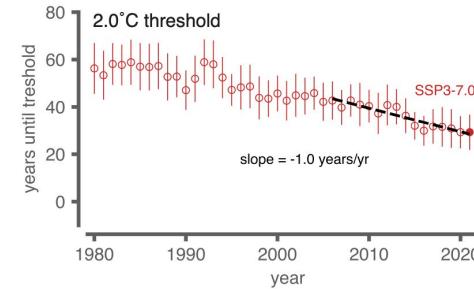
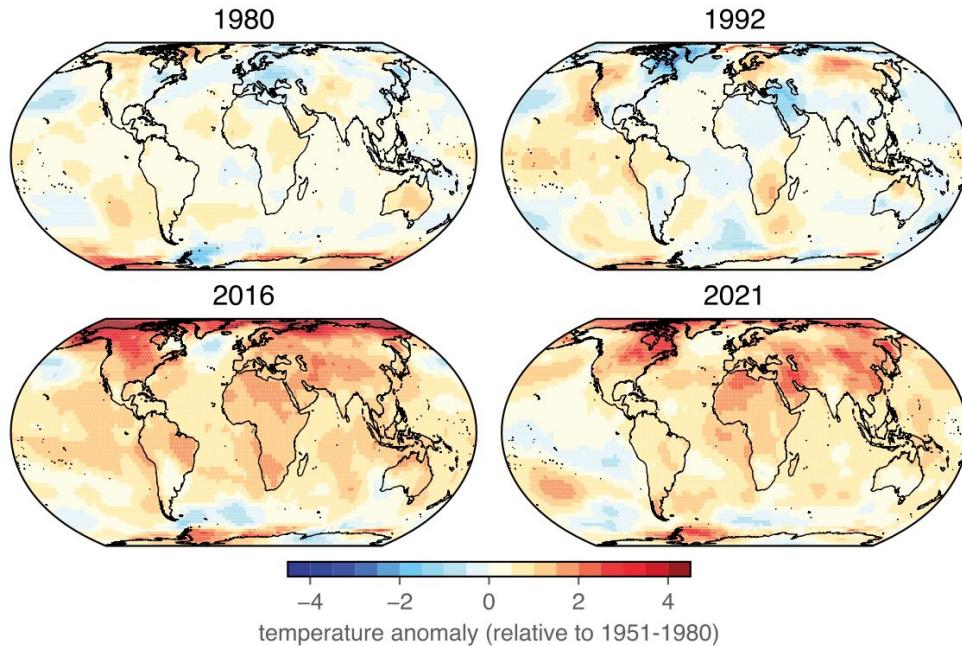
Observations

Berkeley Earth Surface Temperature

**Use the trained AI model to predict thresholds based
on maps of the observed climate**



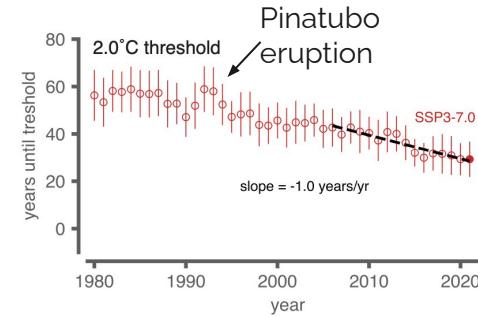
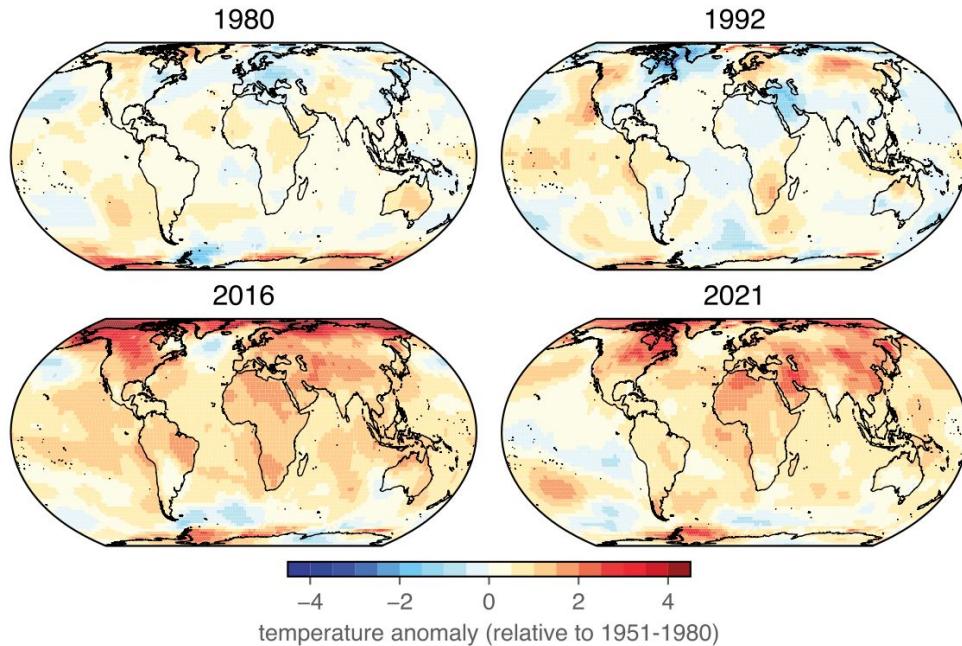
Diffenbaugh & Barnes (2023)



Use the trained AI model to predict thresholds based on maps of the observed climate



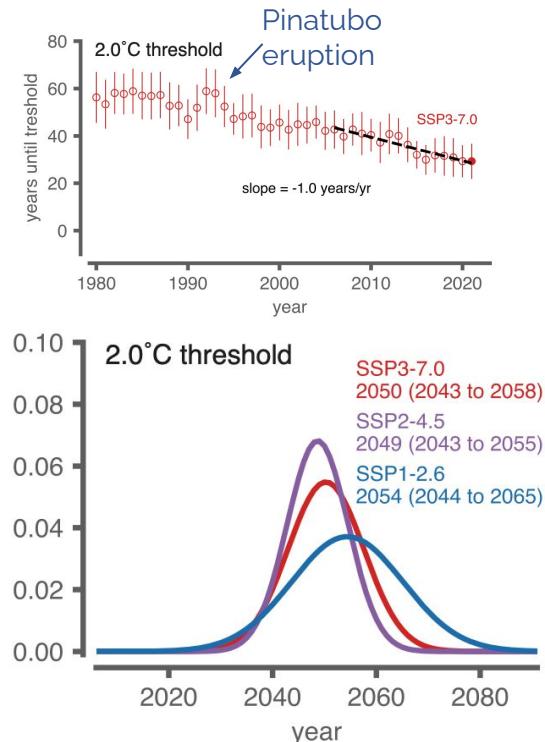
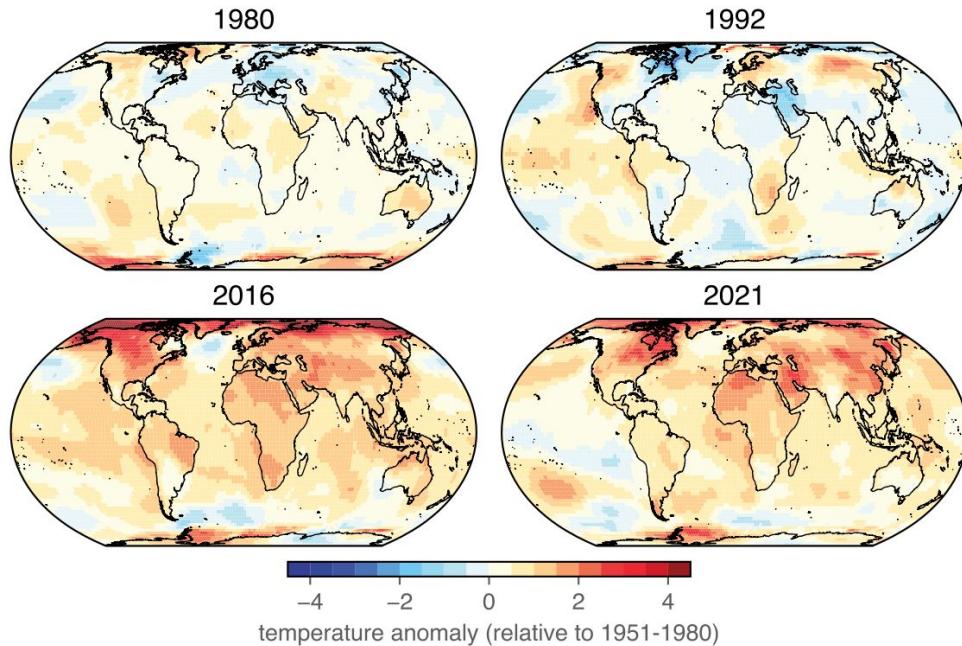
Diffenbaugh & Barnes (2023)



Use the trained AI model to predict thresholds based on maps of the observed climate



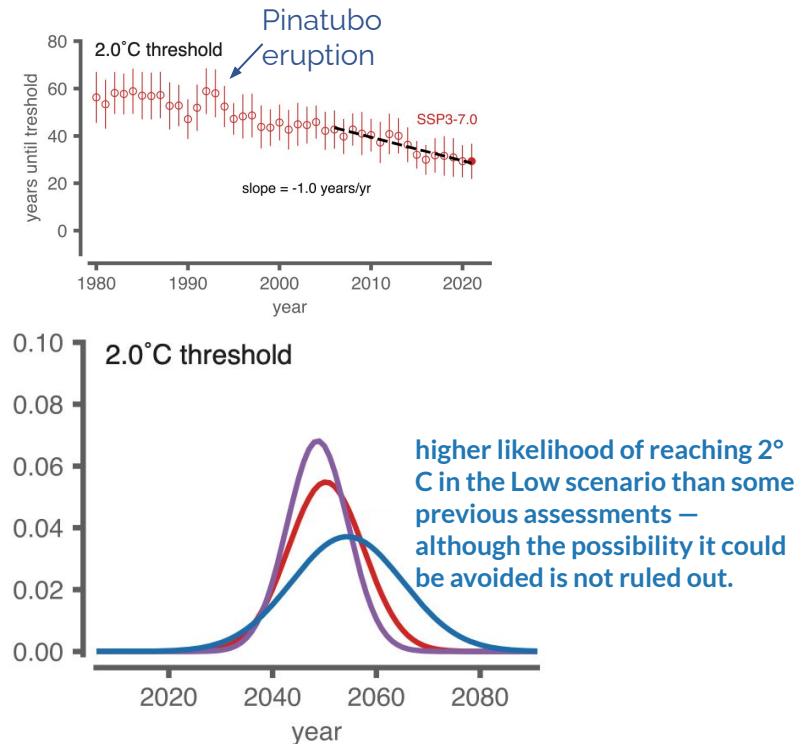
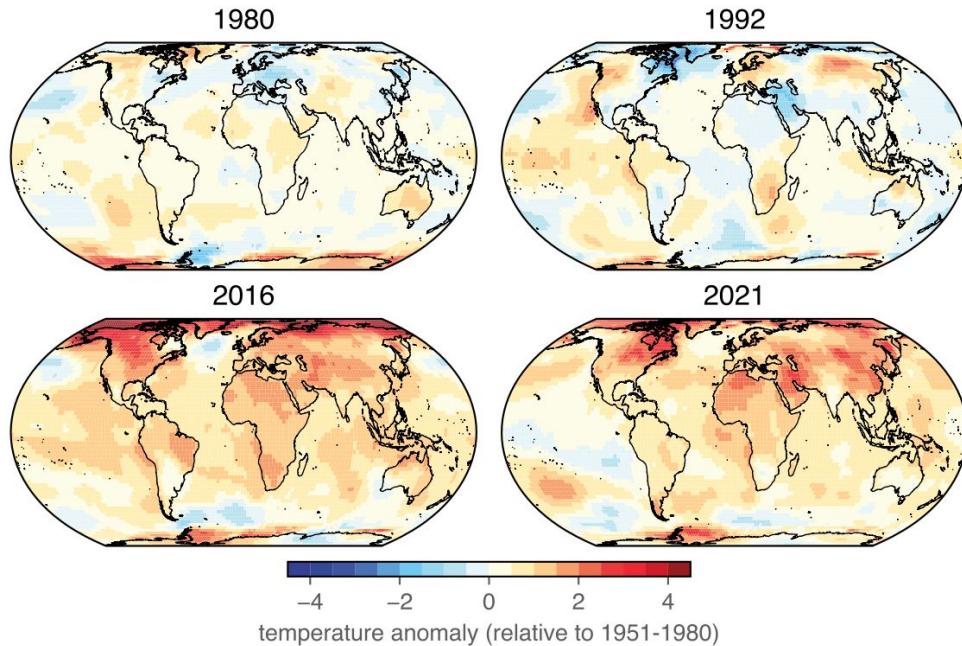
Diffenbaugh & Barnes (2023)



Use the trained AI model to predict thresholds based on maps of the observed climate



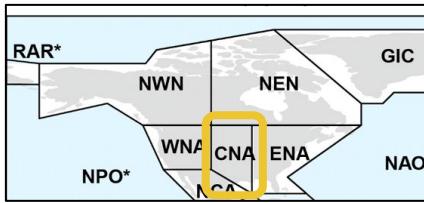
Diffenbaugh & Barnes (2023)



Use the trained AI model to predict thresholds based on maps of the observed climate



Diffenbaugh & Barnes (2023)



Transfer Learning

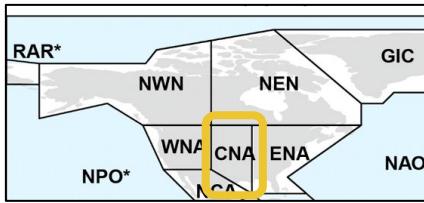
train on climate
model data

predict observations

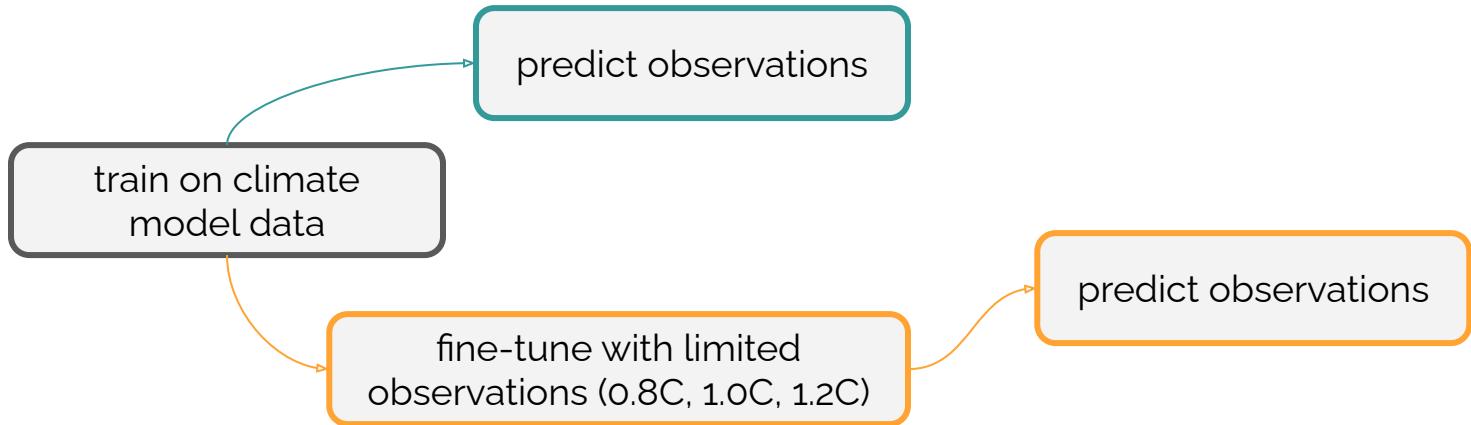
Regional transfer learning provides new insights



Barnes, Diffenbaugh & Seneviratne (2024)



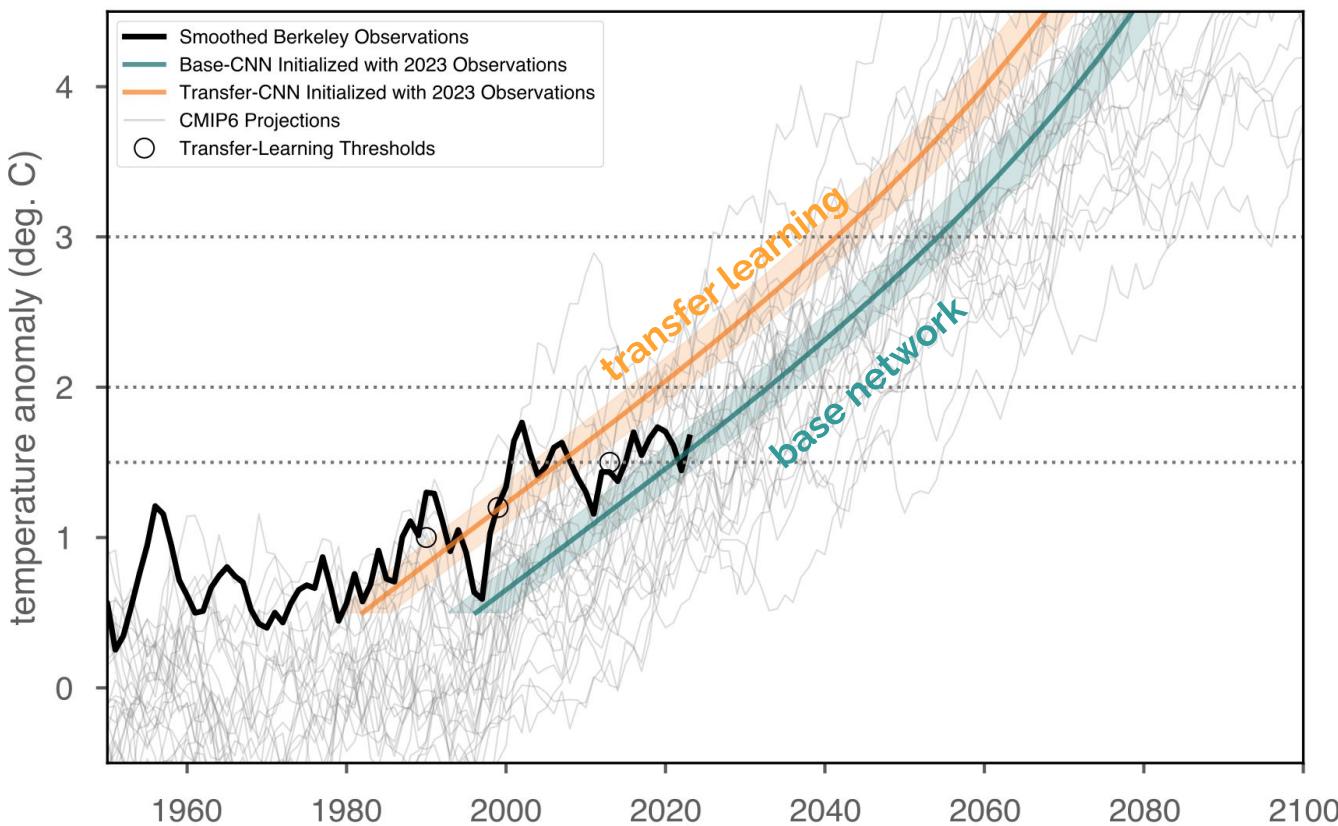
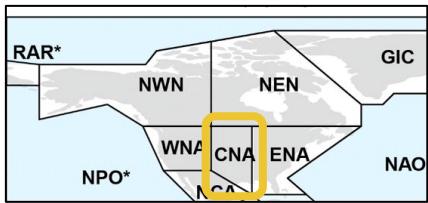
Transfer Learning



Regional transfer learning provides new insights



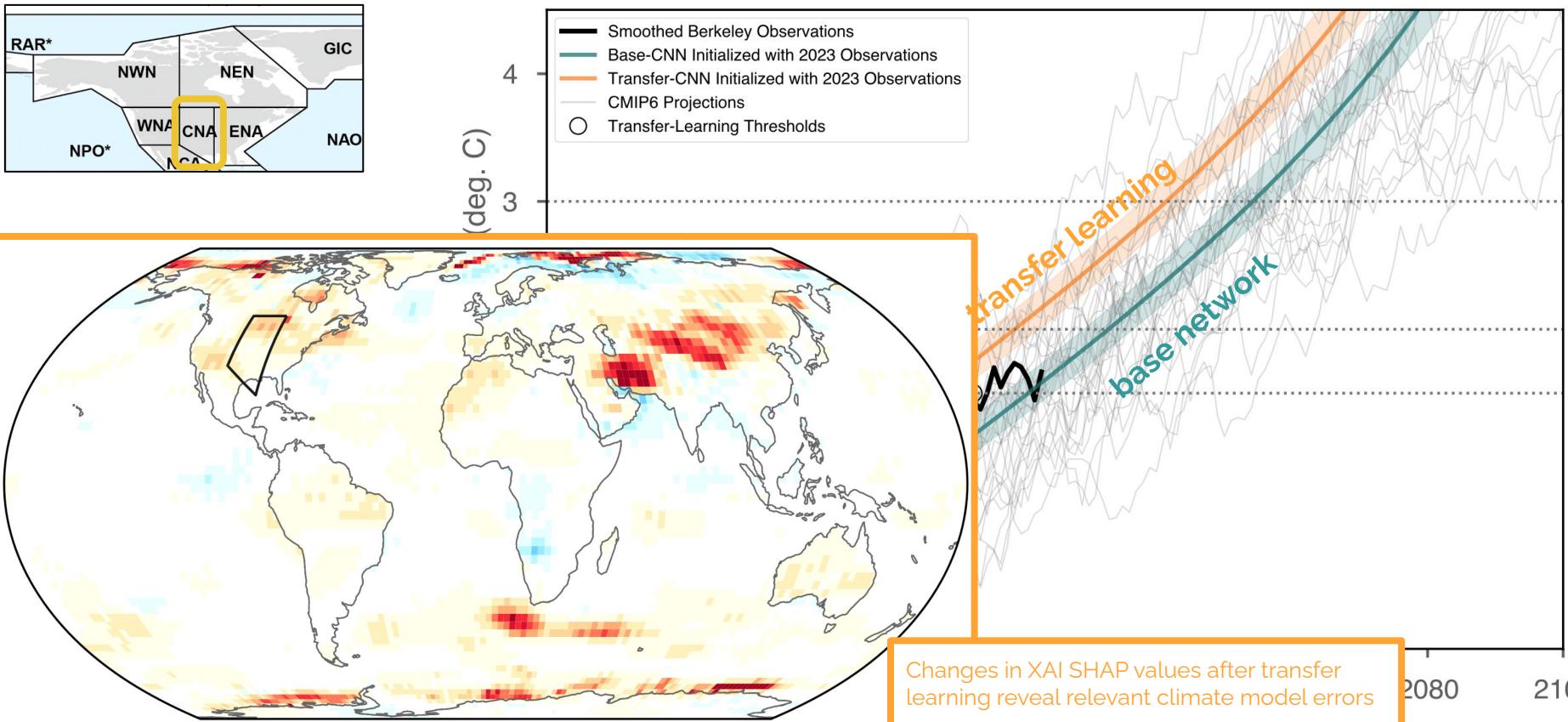
Barnes, Diffenbaugh & Seneviratne (2024)



Regional transfer learning provides new insights



Barnes, Diffenbaugh & Seneviratne (2024)

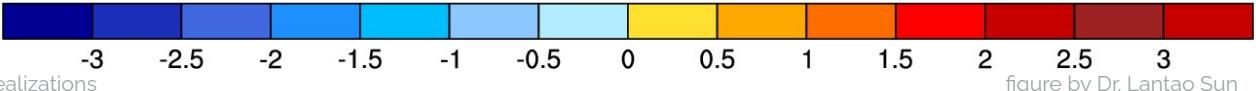
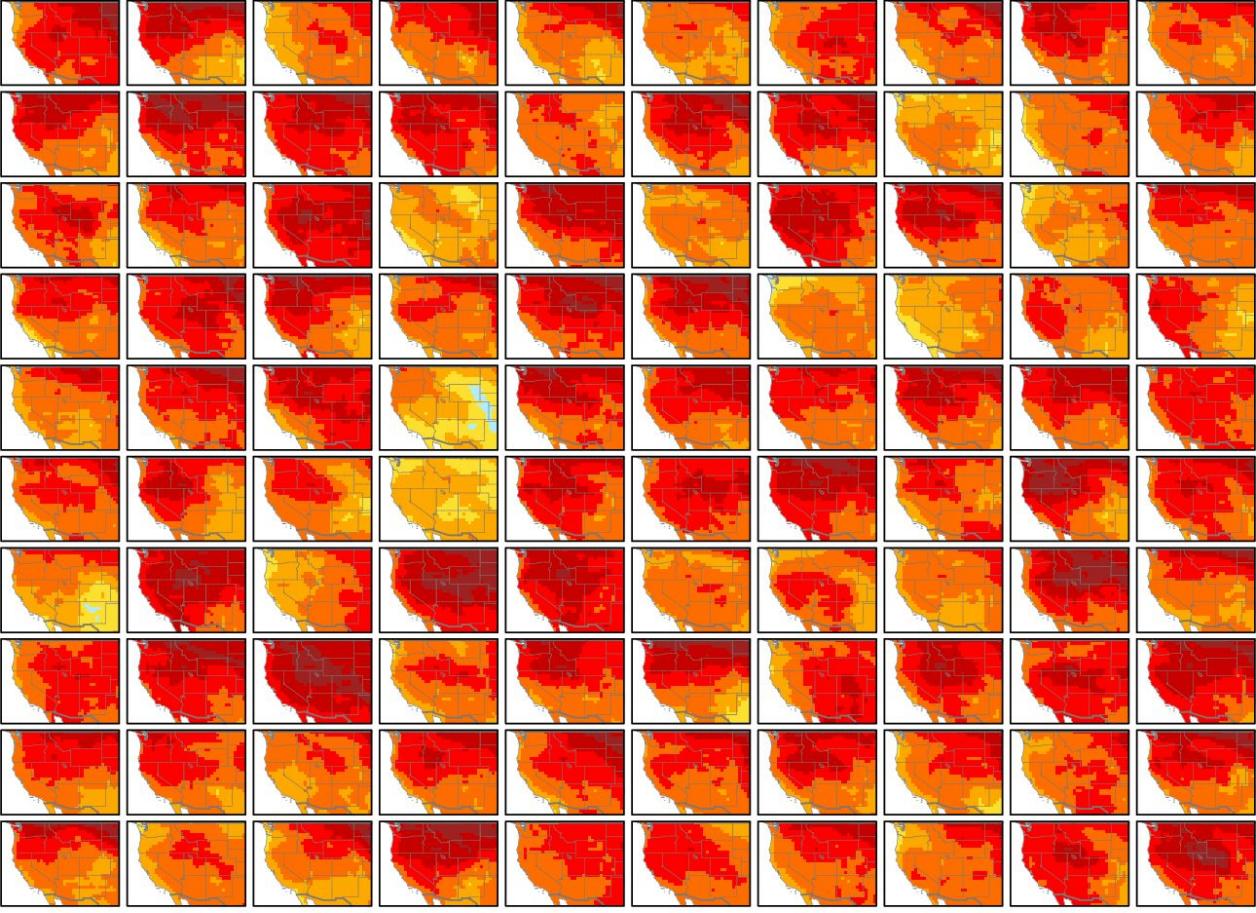


Regional transfer learning provides new insights



Barnes, Diffenbaugh & Seneviratne (2024)

Surface Temperature Linear Trends **2021-2050**

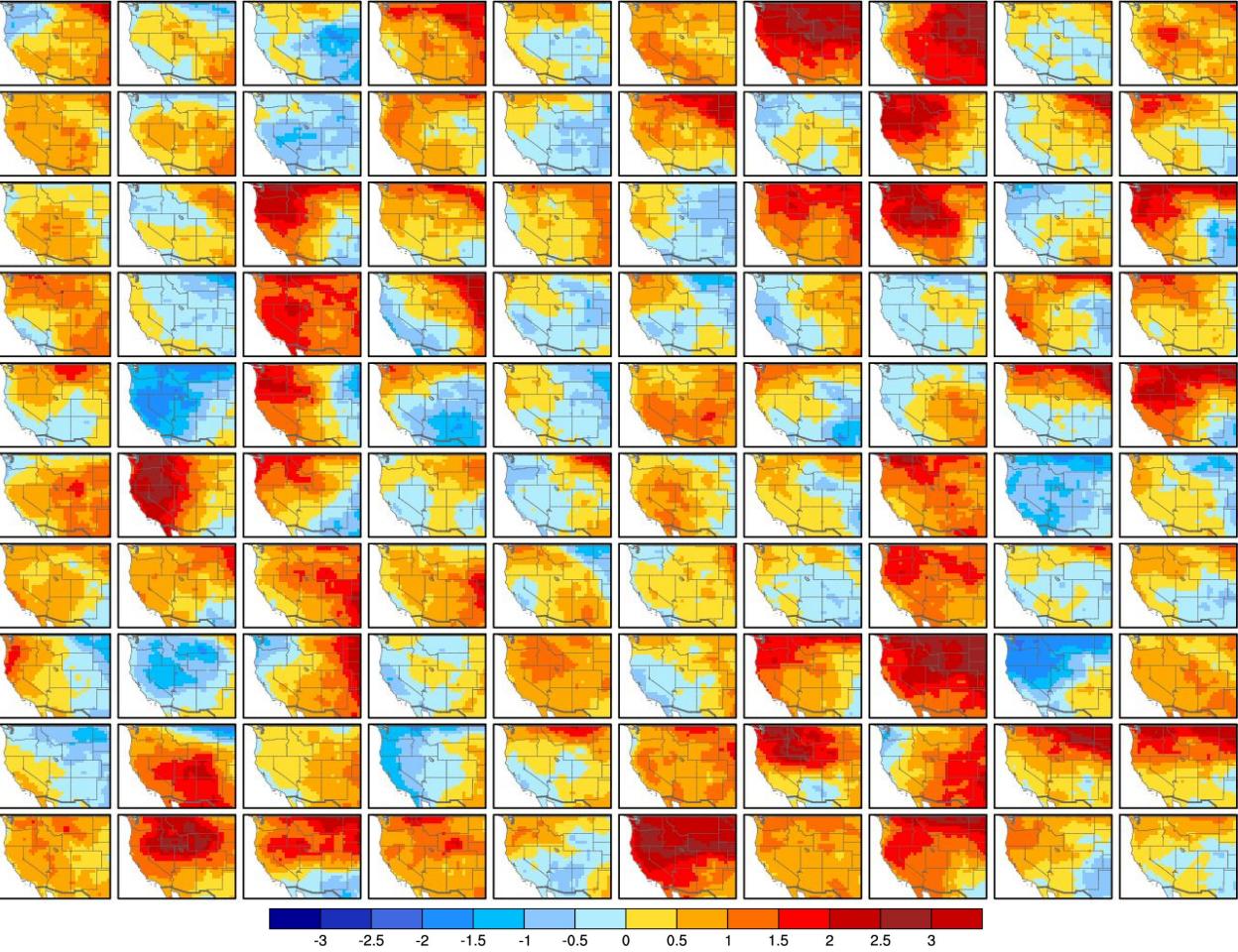


Annual mean surface temperature
As simulated by the CESM2 climate model over 100 different realizations

figure by Dr. Lantao Sun

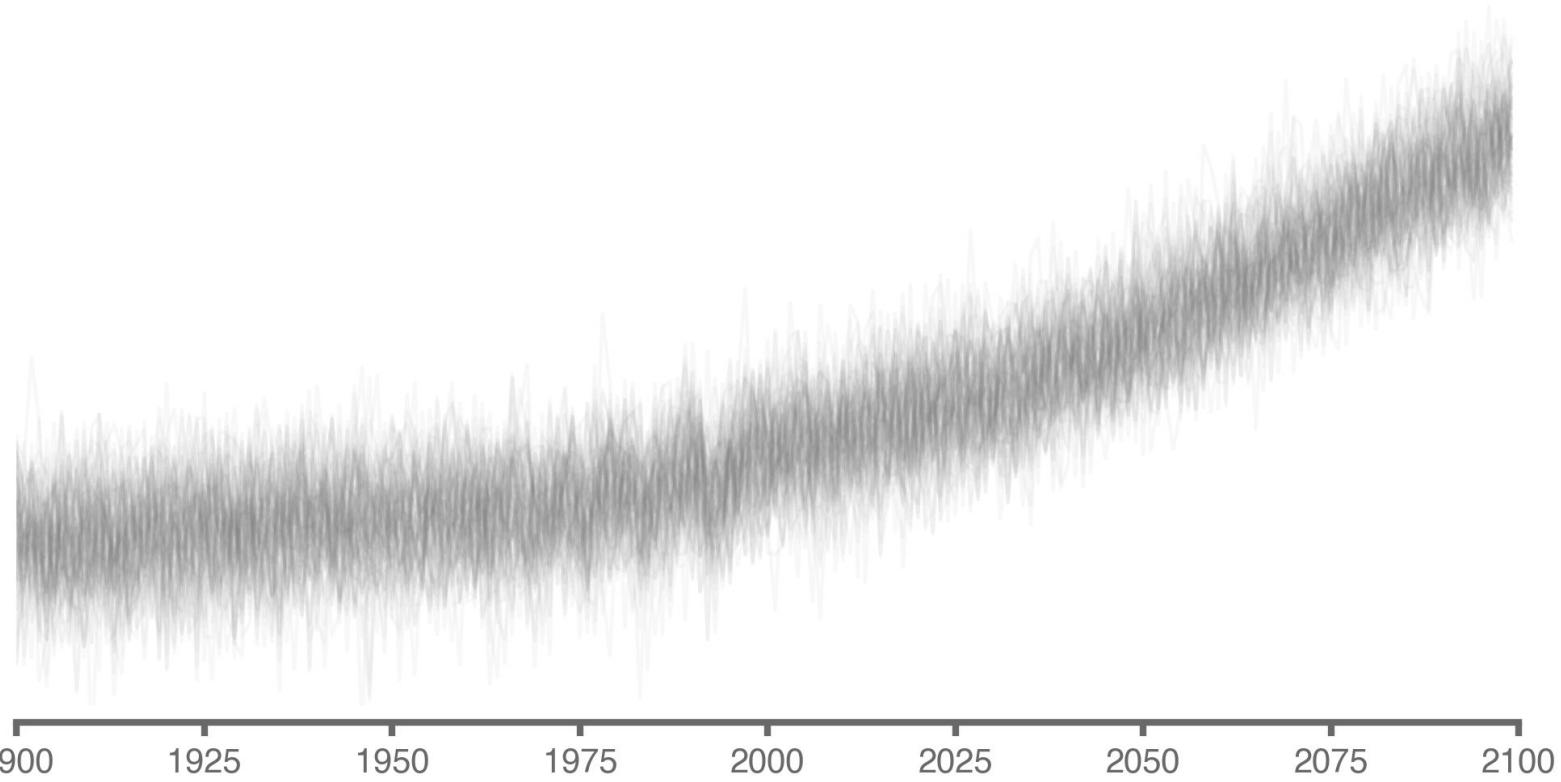
Surface Temperature Linear Trends **2021-2030**

Natural variability is much larger than the climate change signal on regional scales and confounds our knowledge of what the near future holds

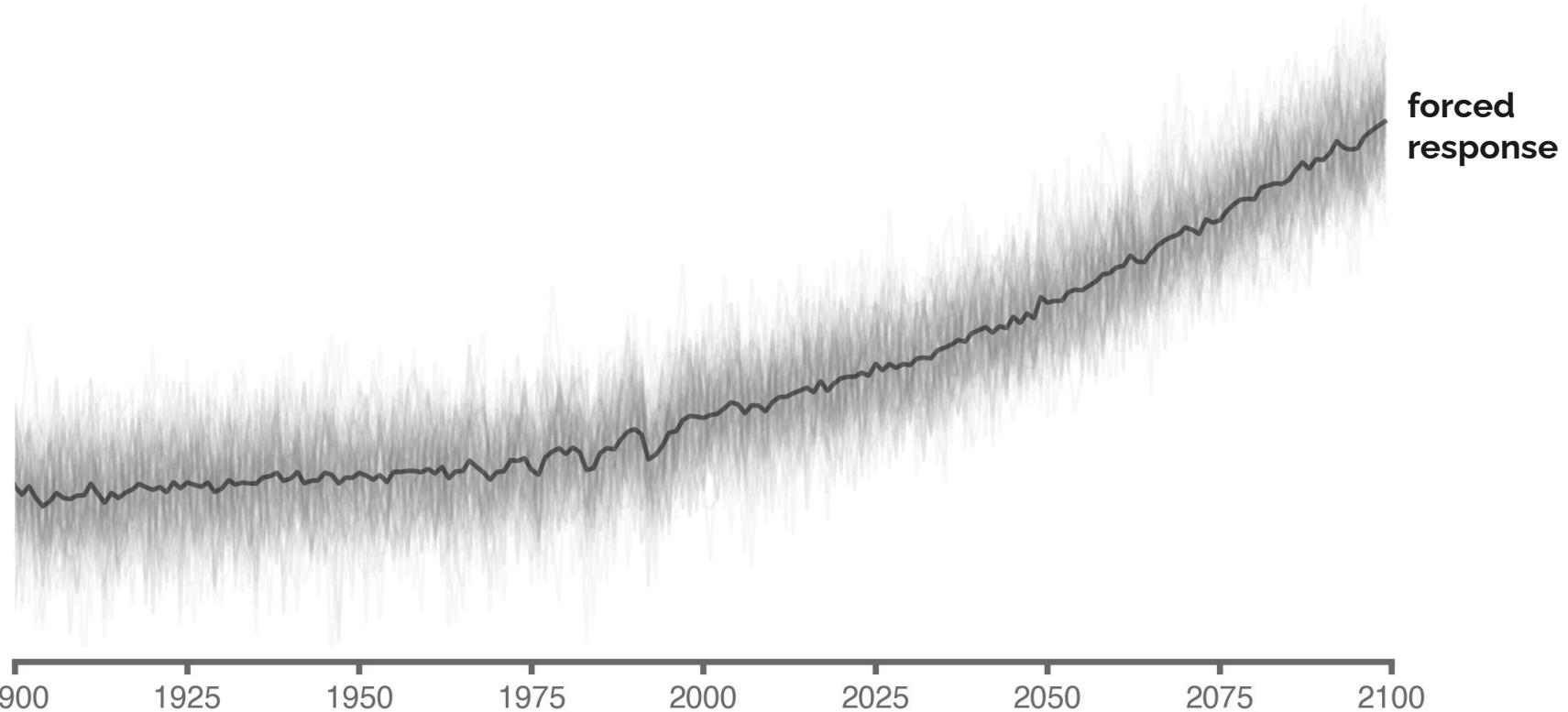


Annual mean surface temperature
As simulated by the CESM2 climate model over 100 different realizations

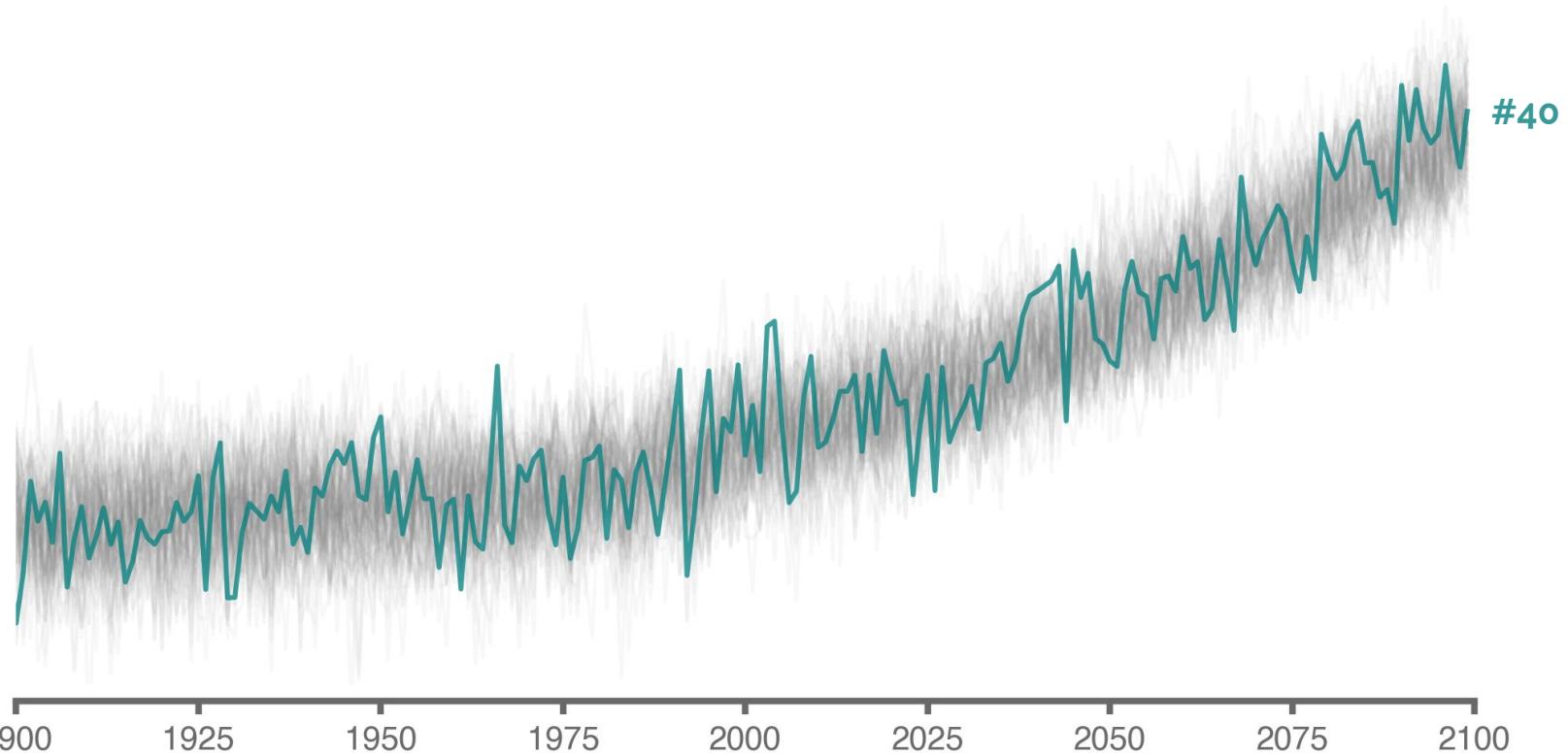
figure by Dr. Lantao Sun



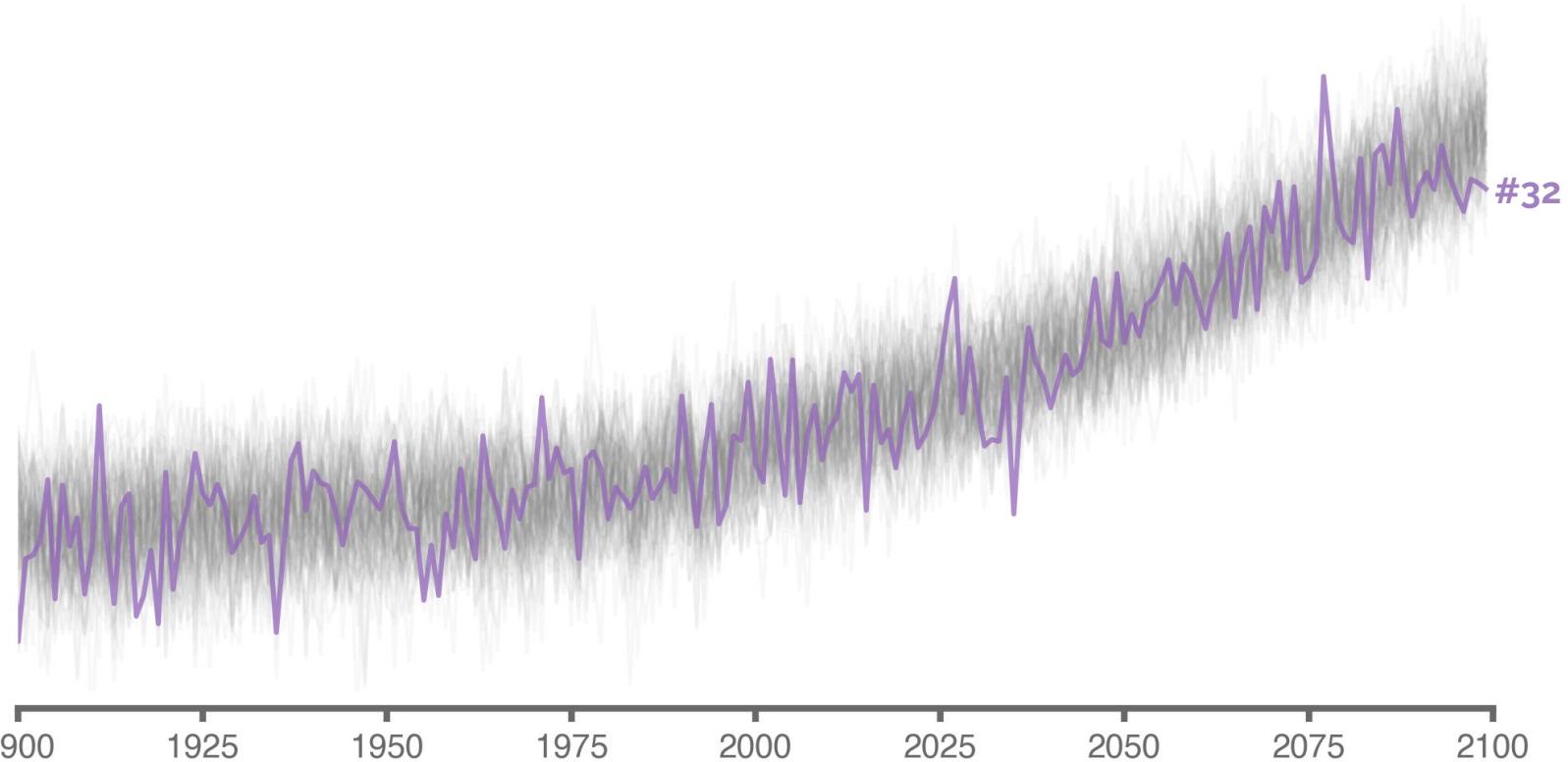
Surface temperature over Chicago, IL
MPI-ESM Large Ensemble; historical + RCP8.5



Surface temperature over Chicago, IL
MPI-ESM Large Ensemble; historical + RCP8.5

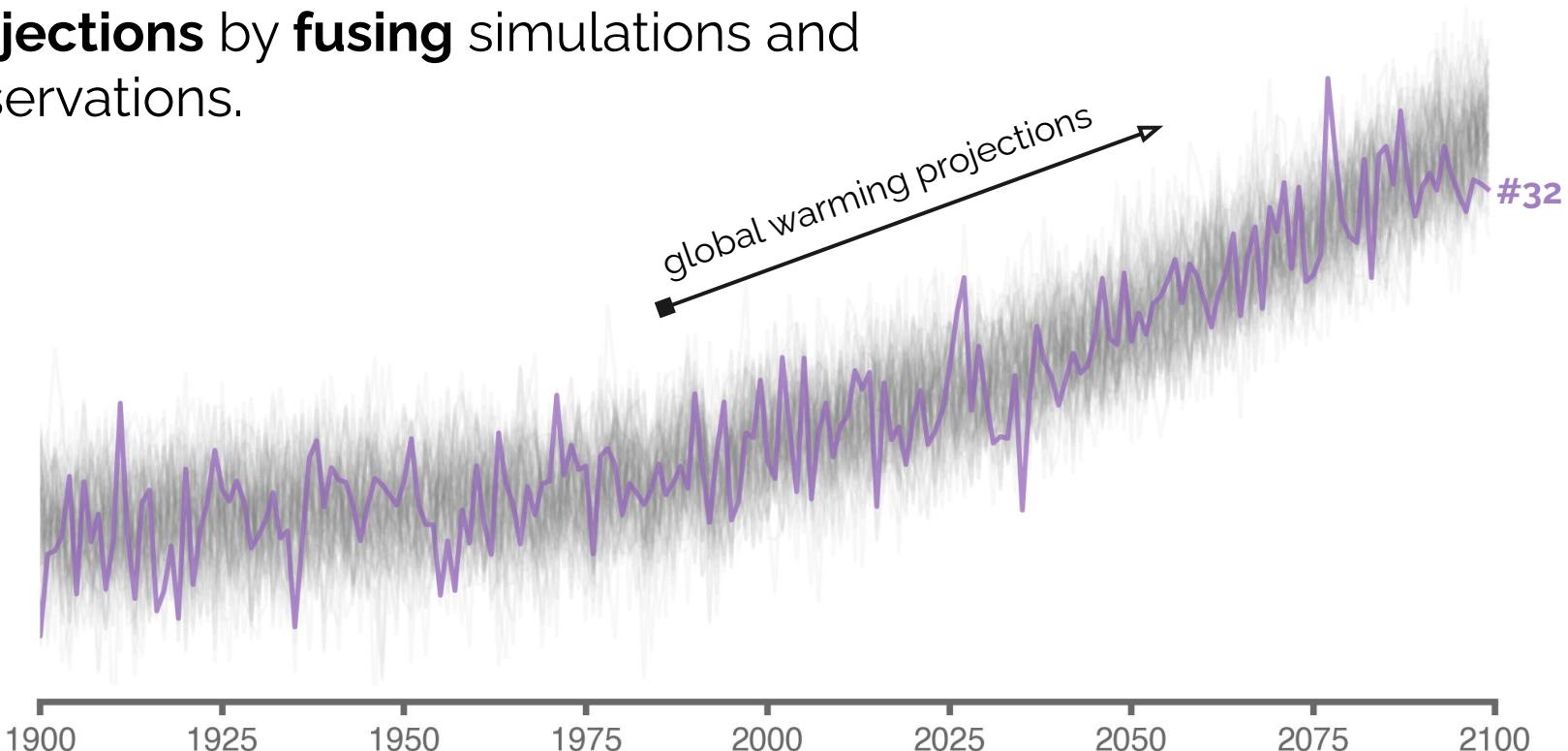


Surface temperature over Chicago, IL
MPI-ESM Large Ensemble; historical + RCP8.5



Surface temperature over Chicago, IL
MPI-ESM Large Ensemble; historical + RCP8.5

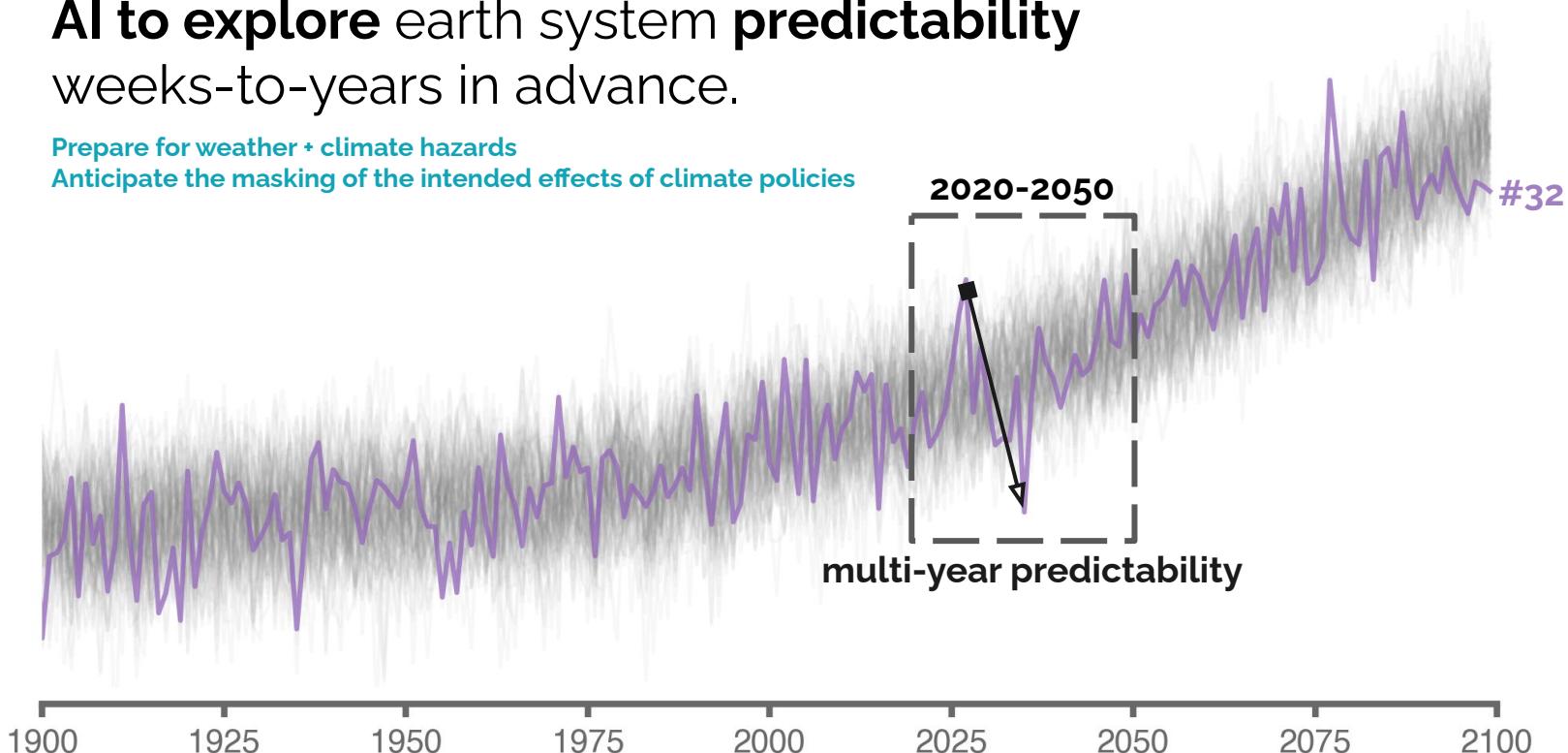
AI to leverage imperfect climate models to better **constrain future projections** by **fusing** simulations and observations.



AI to explore earth system **predictability** weeks-to-years in advance.

Prepare for weather + climate hazards

Anticipate the masking of the intended effects of climate policies



1900

1925

1950

1975

2000

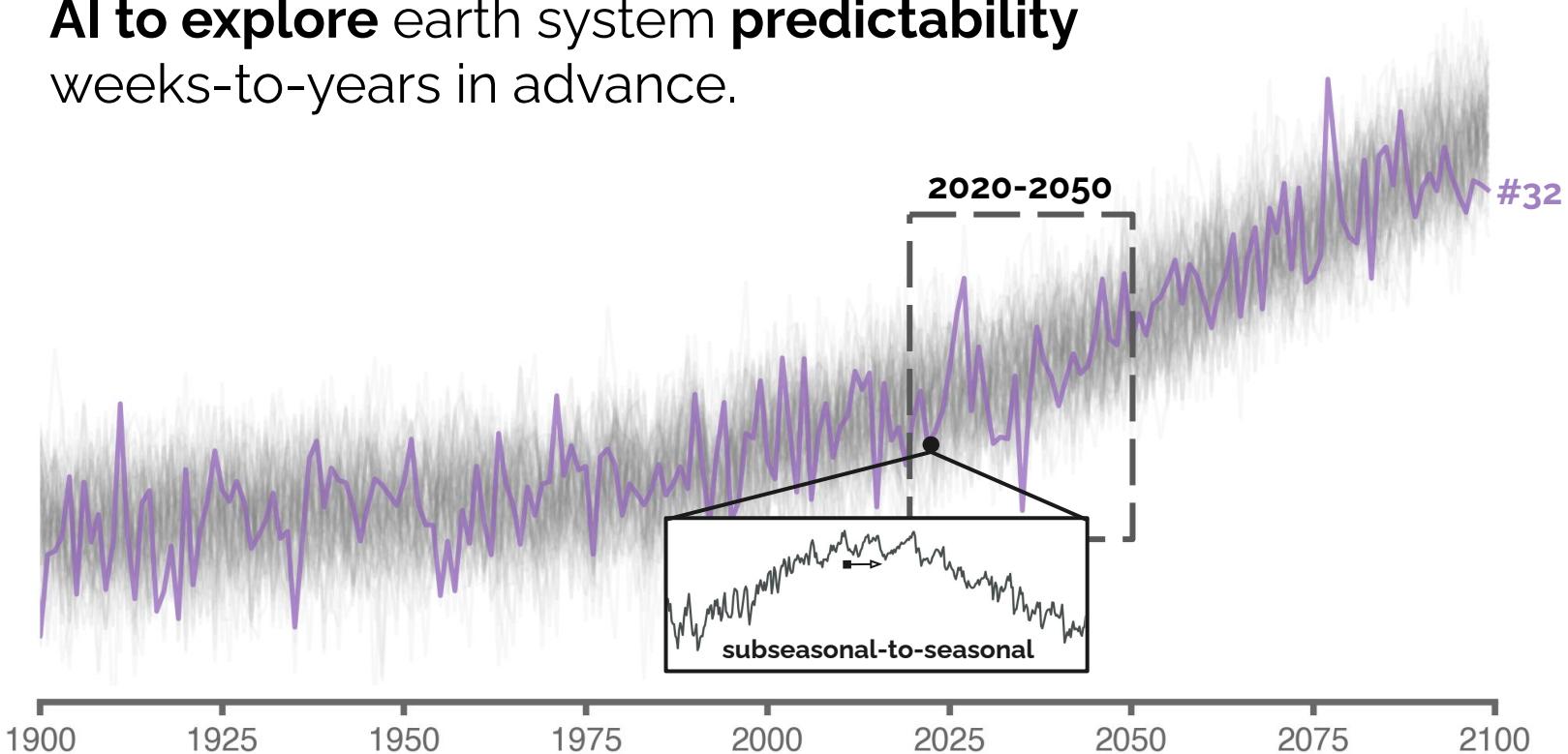
2025

2050

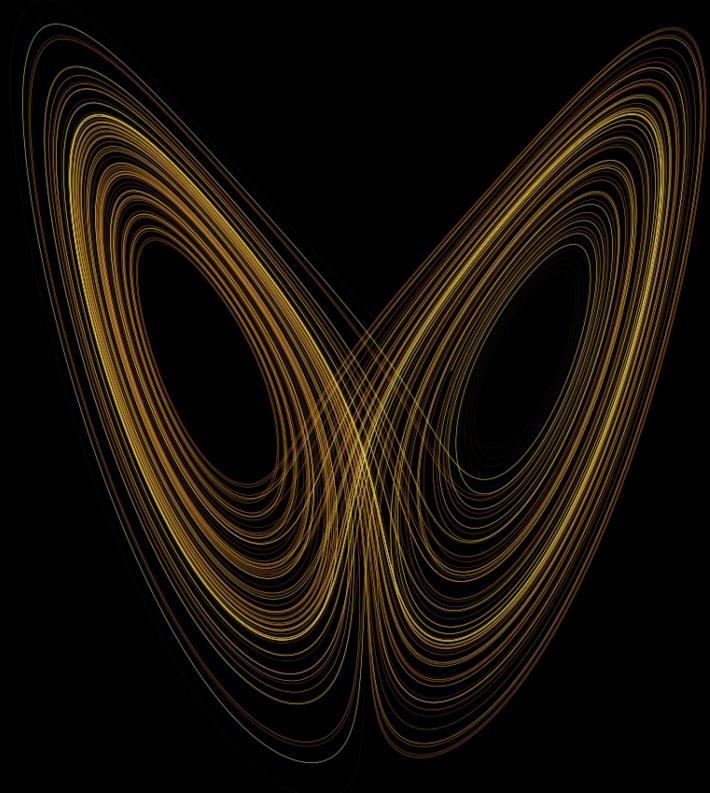
2075

2100

AI to explore earth system **predictability**
weeks-to-years in advance.

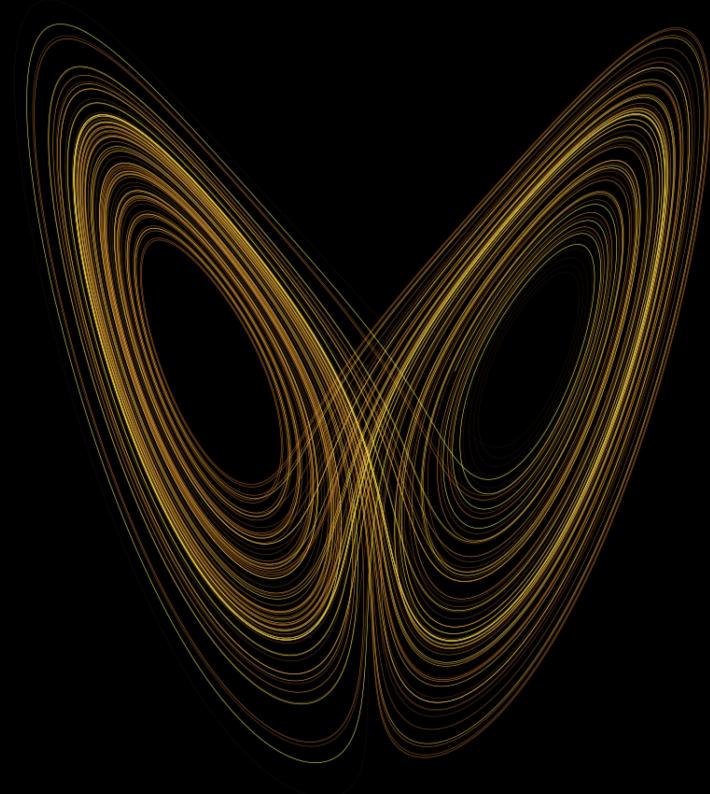


**But, climate prediction is
incredibly challenging.
We cannot expect to
make perfect predictions
all of the time.**



**But, climate prediction is
incredibly challenging.
We cannot expect to
make perfect predictions
all of the time.**

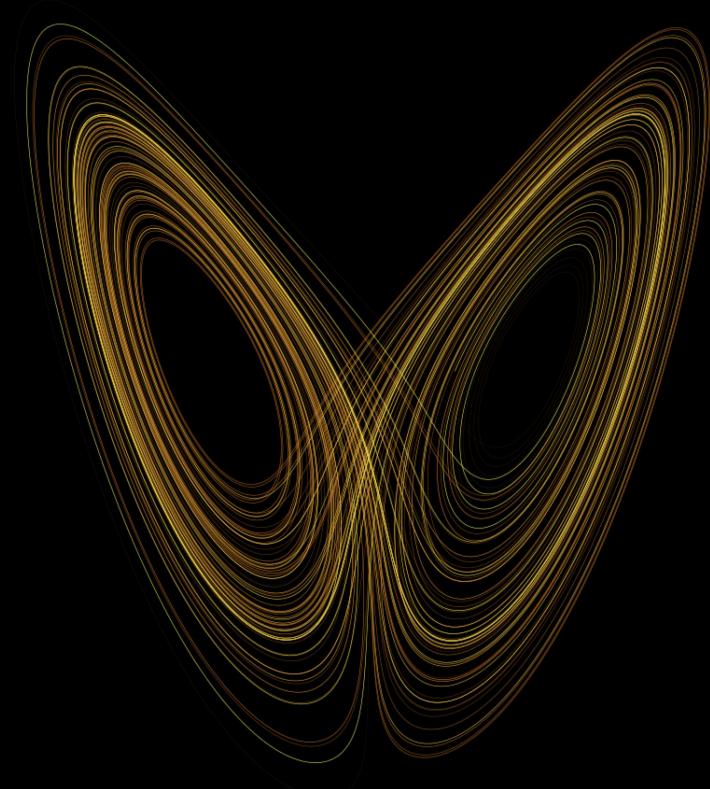
Instead, we must look for specific states of the earth system, i.e. "forecasts of opportunity", that lead to enhanced predictable behavior.



**But, climate prediction is
incredibly challenging.
We cannot expect to
make perfect predictions
all of the time.**

Instead, we must look for specific states of the earth system, i.e. "forecasts of opportunity", that lead to enhanced predictable behavior.

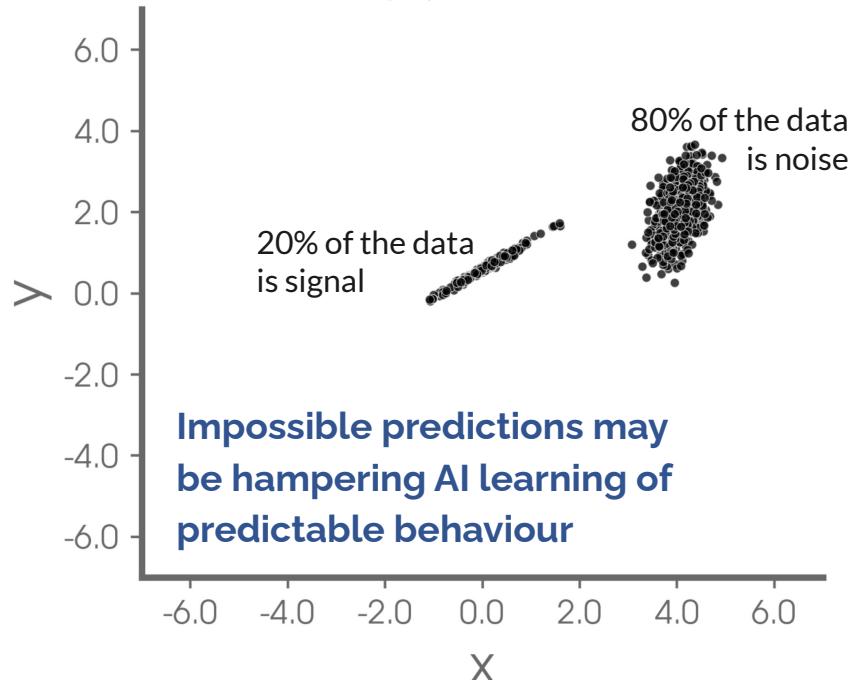
AI can help us with this.



**But, climate prediction is
incredibly challenging.
We cannot expect to
make perfect predictions
all of the time.**

Instead, we must look for specific states of the earth system, i.e. "forecasts of opportunity", that lead to enhanced predictable behavior.

AI can help us with this.

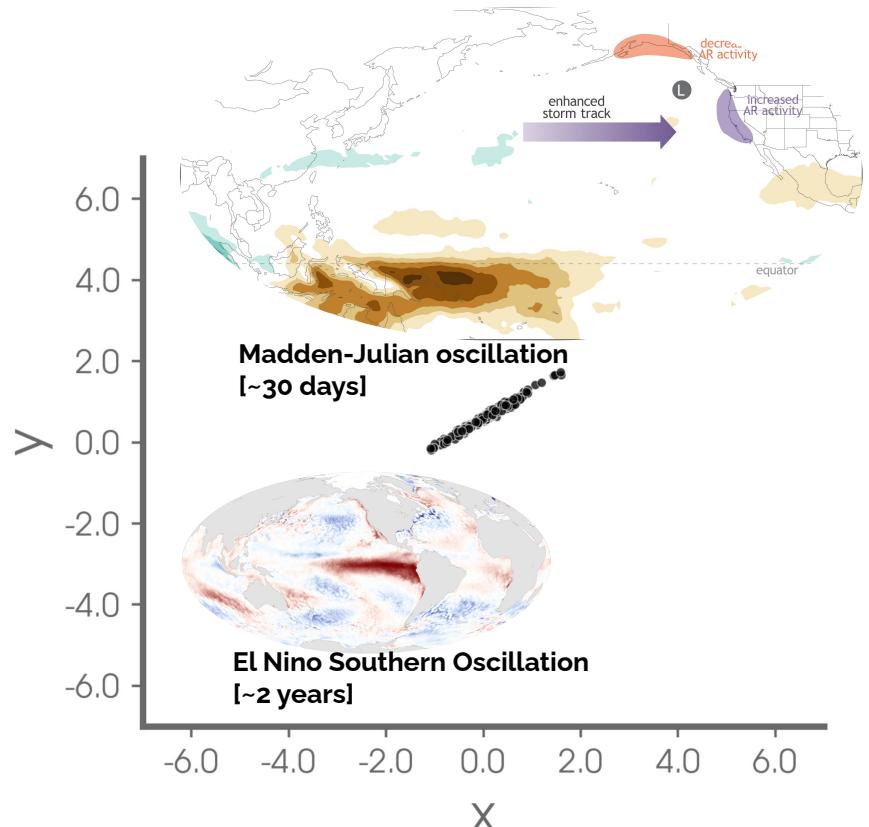


Barnes, Barnes and Gordillo (2021)
Barnes and Barnes (2021a, 2021b)

**But, climate prediction is
incredibly challenging.
We cannot expect to
make perfect predictions
all of the time.**

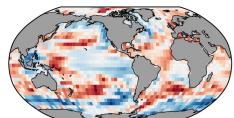
Instead, we must look for specific states
of the earth system, i.e. "forecasts of
opportunity", that lead to enhanced
predictable behavior.

AI can help us with this.

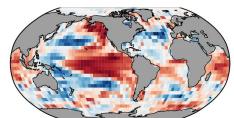


Barnes, Barnes and Gordillo (2021)
Barnes and Barnes (2021a, 2021b)

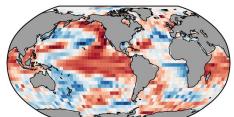
past sea-surface temperatures



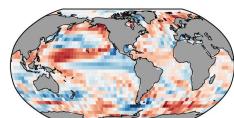
3-8 years before



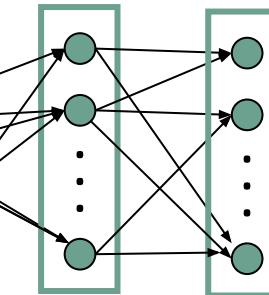
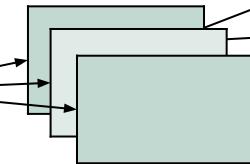
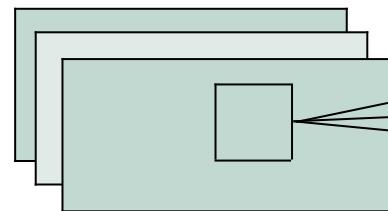
2-3 years before



1-2 years before



0-1 years before



future sea surface temperatures*
for one grid point
[0-5 years]

warm

neutral

cool

*can predict a range of variables

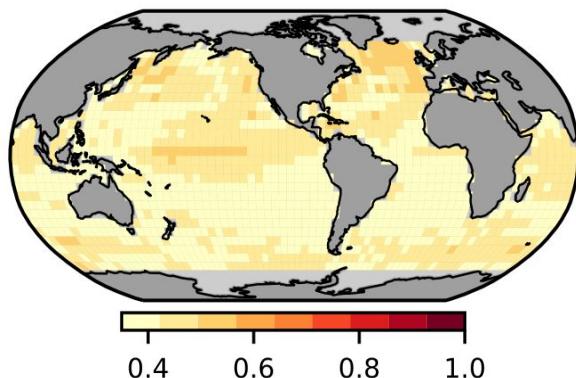
Predict ocean temperatures 5 years later



Davenport, Barnes & Gordon (2024)

CLIMATE MODEL DATA

Overall Accuracy



Trained on climate model **MPI-ESM-1-2-LR [3,630 years of data]**
Evaluated on climate model **MPI-ESM-1-2-LR**

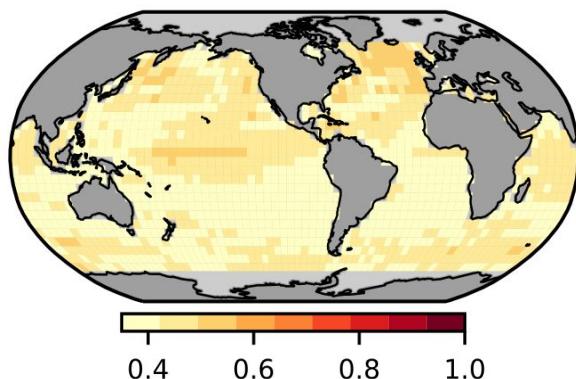
Focusing on when the AI is most confident leads to accurate predictions



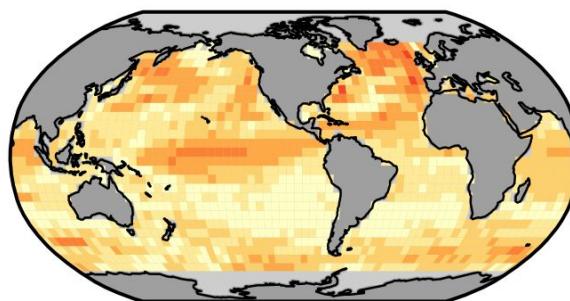
Davenport, Barnes & Gordon (2024)

CLIMATE MODEL DATA

Overall Accuracy



Accuracy for 40% most confident predictions



Trained on climate model **MPI-ESM-1-2-LR [3,630 years of data]**

Evaluated on climate model **MPI-ESM-1-2-LR**

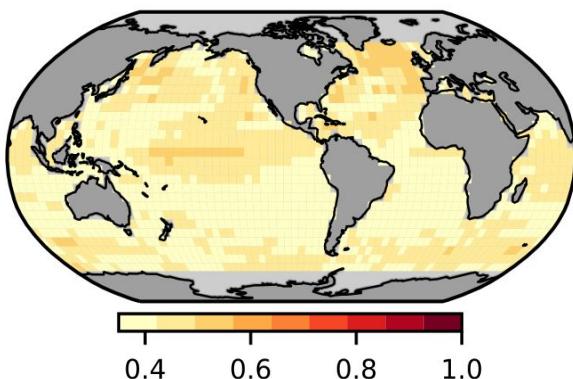
Focusing on when the AI is most confident leads to accurate predictions



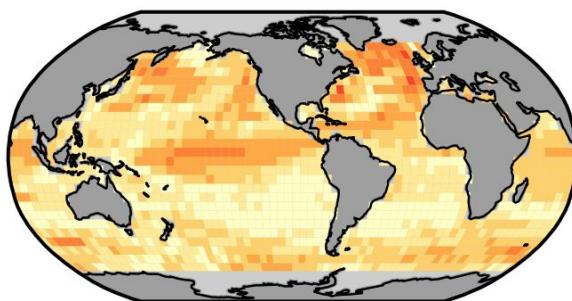
Davenport, Barnes & Gordon (2024)

CLIMATE MODEL DATA

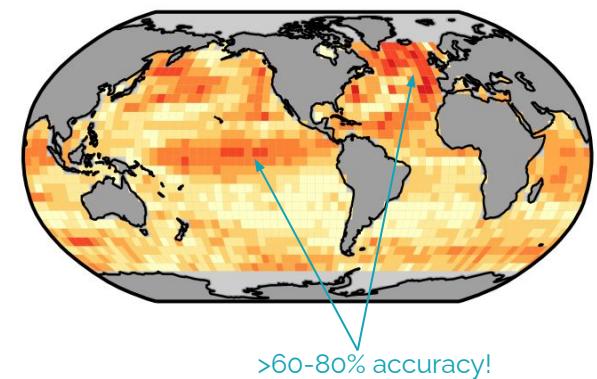
Overall Accuracy



Accuracy for 40% most confident predictions



Accuracy for 20% most confident predictions



Trained on climate model **MPI-ESM-1-2-LR [3,630 years of data]**
Evaluated on climate model **MPI-ESM-1-2-LR**

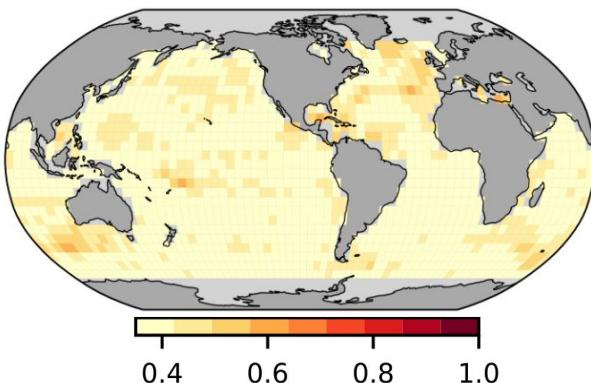
Focusing on when the AI is most confident leads to accurate predictions



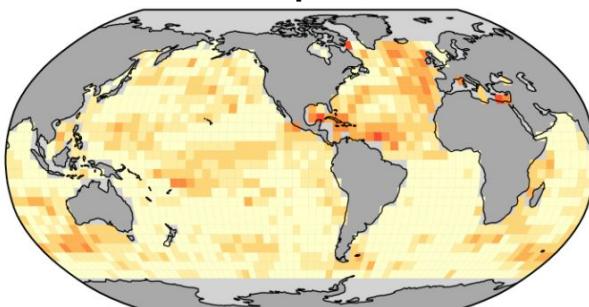
Davenport, Barnes & Gordon (2024)

OBSERVATIONS

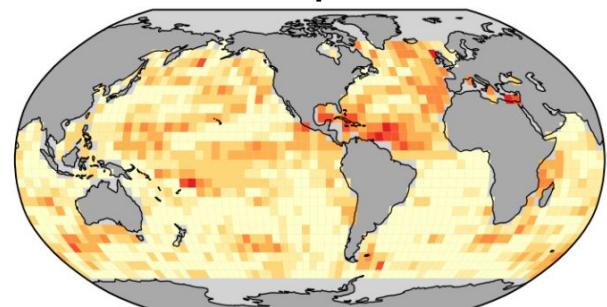
Overall Accuracy



Accuracy for 40% most confident predictions



Accuracy for 20% most confident predictions



Trained on climate model **MPI-ESM-1-2-LR** [3,630 years of data]
Evaluated on **observations** [ERSSTv5; 169 years of data]

Leveraging climate model data provides accurate predictions of the real world



Davenport, Barnes & Gordon (2024)

But how is the network making its prediction?

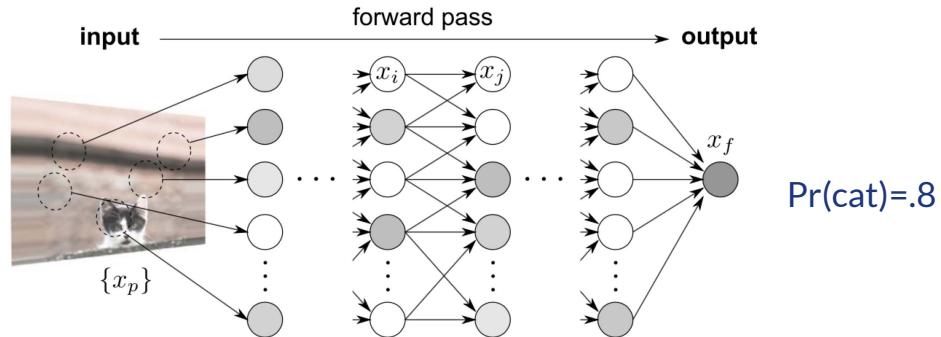
...what if we could learn which regions of the globe are most relevant to each prediction?



XAI Attribution Methods

Attribution heatmaps are largely consistent with how many climate scientists pose questions

Prediction
of 1 sample

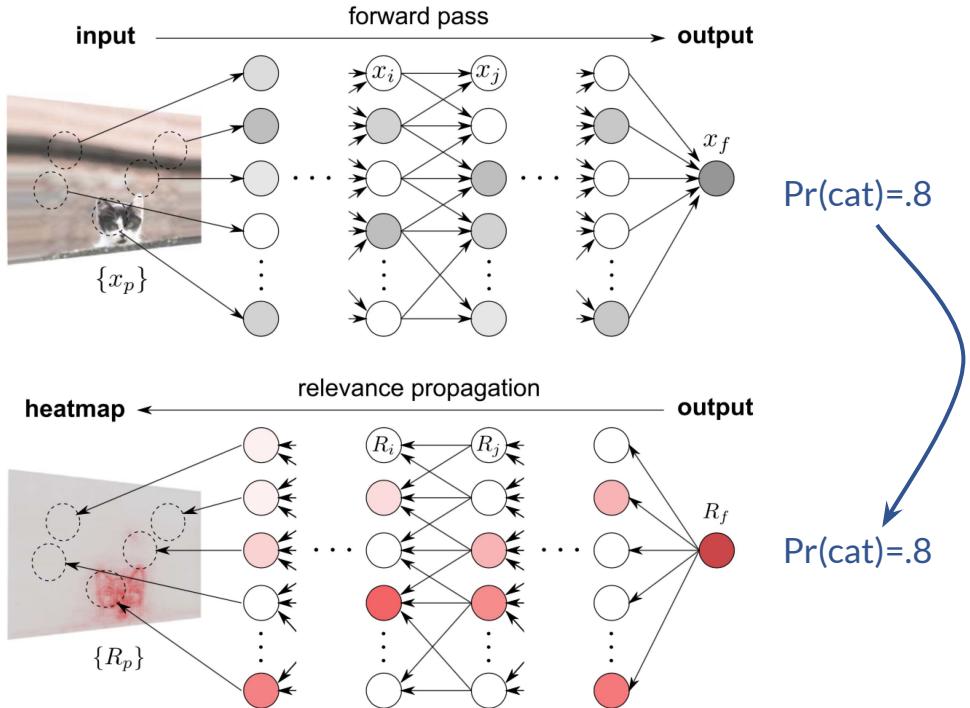


XAI Attribution Methods

Attribution heatmaps are largely consistent with how many climate scientists pose questions

Prediction
of 1 sample

Attribution
of 1 sample



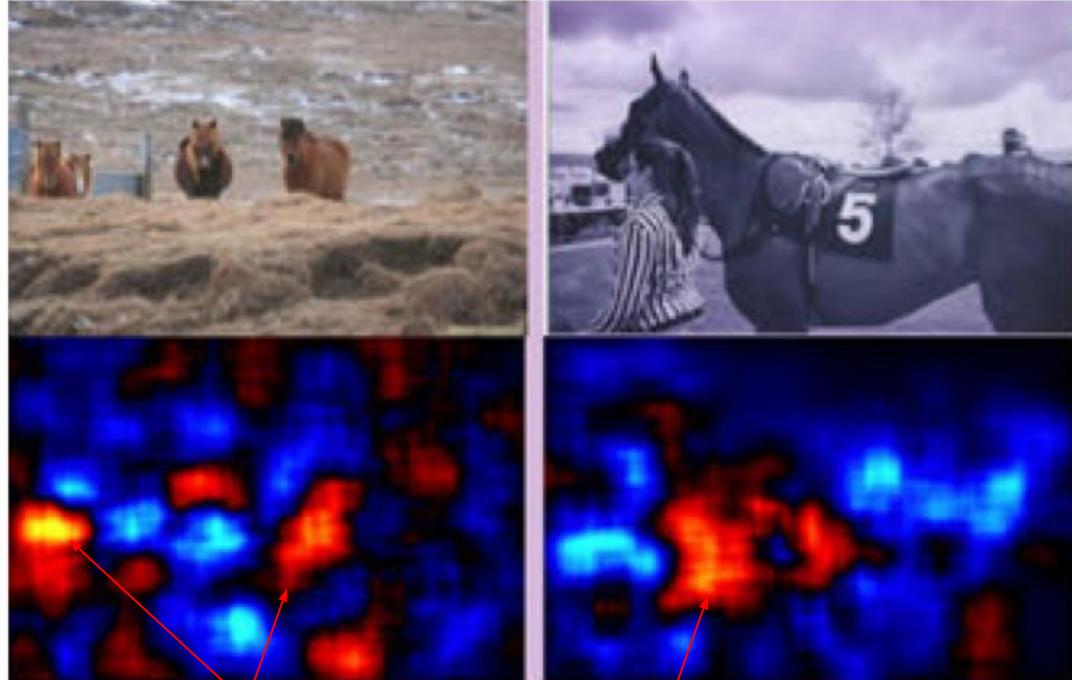
XAI to identify problematic strategies

Is the AI approach predicting the right answers for the right reasons?



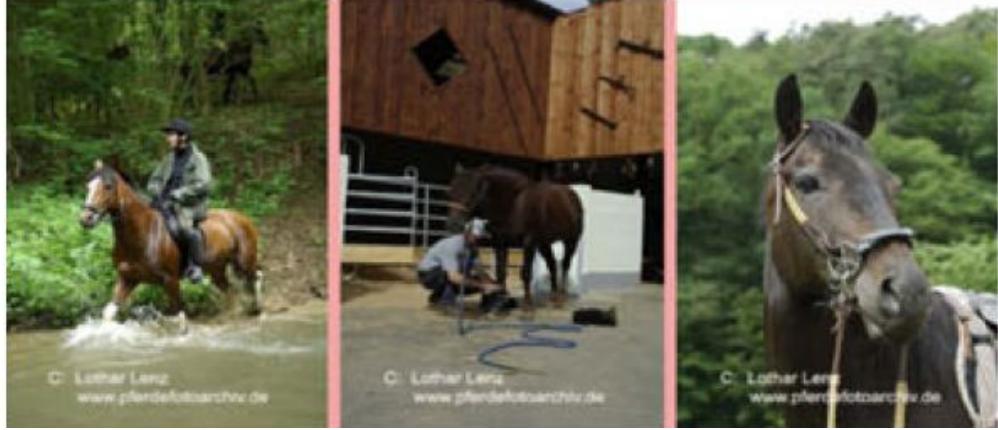
XAI to identify problematic strategies

Is the AI approach predicting the right answers for the right reasons?



XAI to identify problematic strategies

Is the AI approach predicting the right answers for the right reasons?



XAI to identify problematic strategies

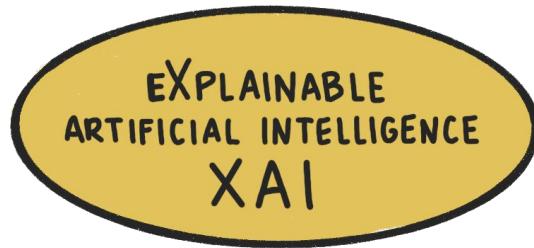
Is the AI approach predicting the right answers for the right reasons?



red shading: relevant regions for making the network think there is a horse

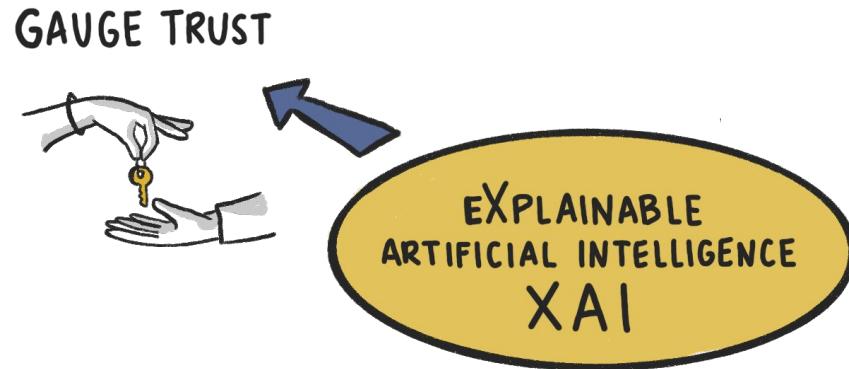
Reasons to care about XAI

A scientist's ultimate goal is typically to understand "why?", but even if you don't care "why?" you should still care about XAI.



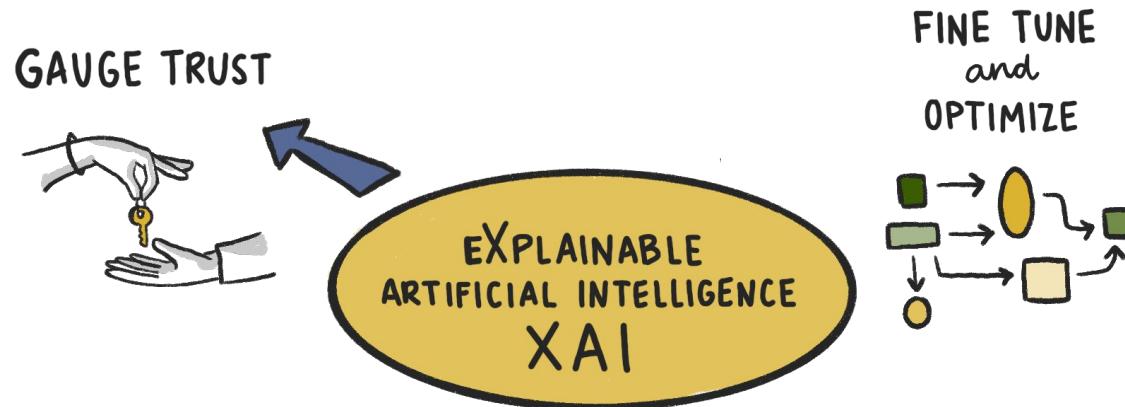
Reasons to care about XAI

A scientist's ultimate goal is typically to understand "why?", but even if you don't care "why?" you should still care about XAI.



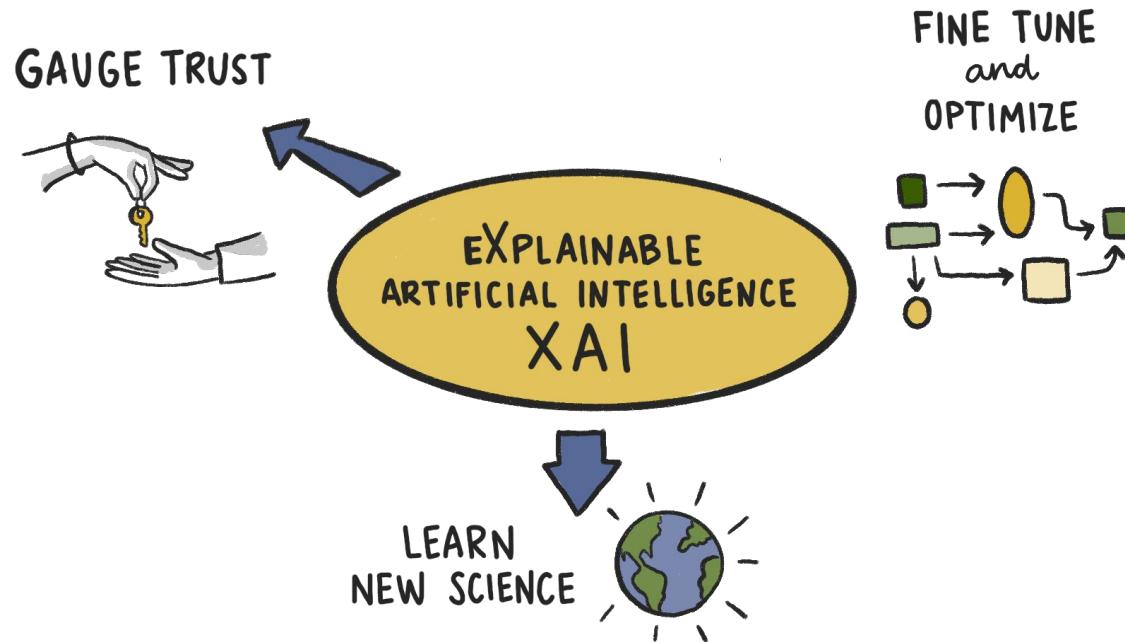
Reasons to care about XAI

A scientist's ultimate goal is typically to understand "why?", but even if you don't care "why?" you should still care about XAI.



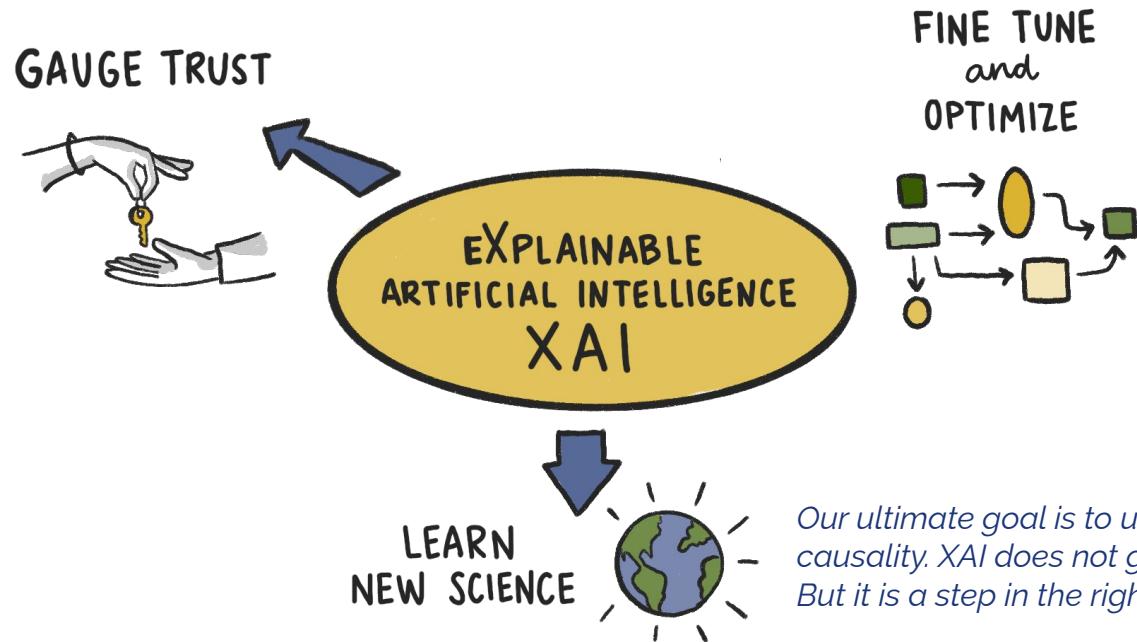
Reasons to care about XAI

A scientist's ultimate goal is typically to understand "why?", but even if you don't care "why?" you should still care about XAI.

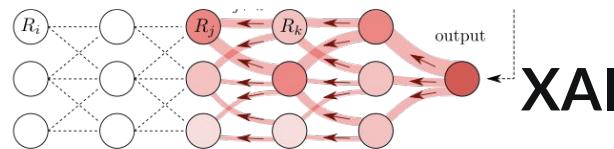
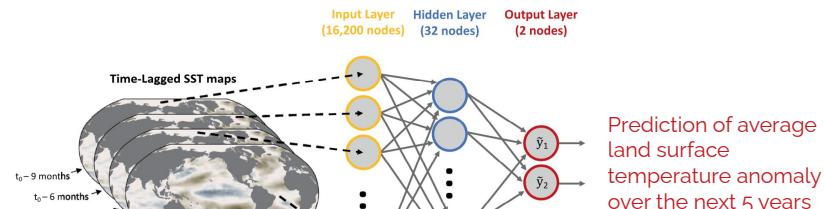


Reasons to care about XAI

A scientist's ultimate goal is typically to understand "why?", but even if you don't care "why?" you should still care about XAI.



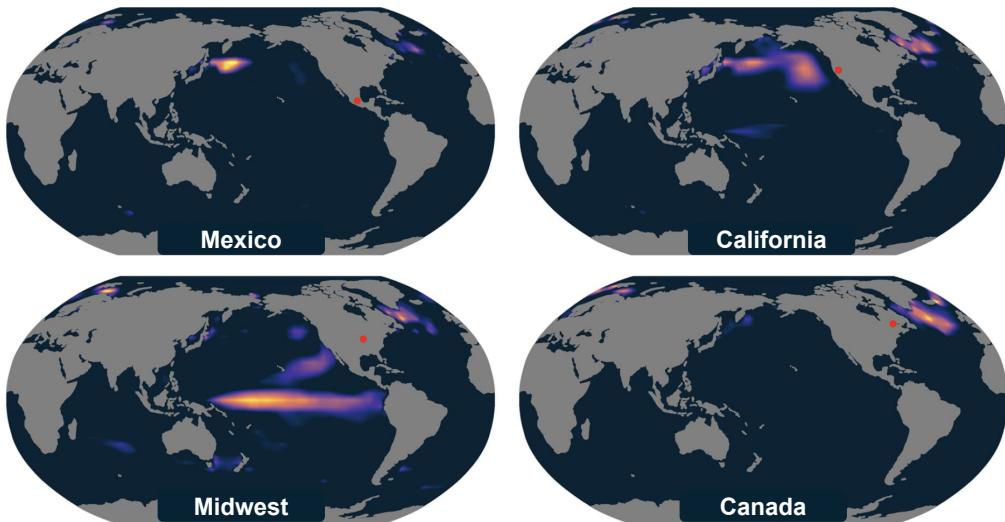
Predicting 5-year average surface temperature at each grid point
Applied to 1200 years of CESM2 control simulation



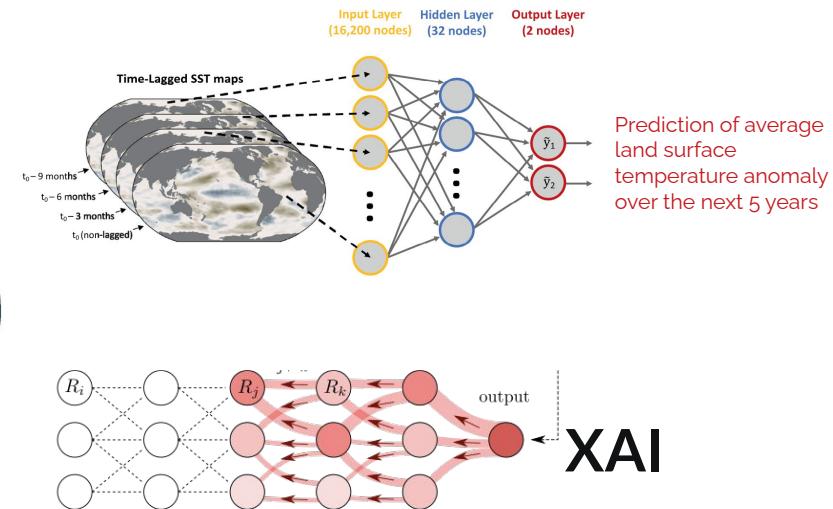
XAI reveals sources of predictability that vary in time and space



Toms, Barnes & Hurrell (2021)



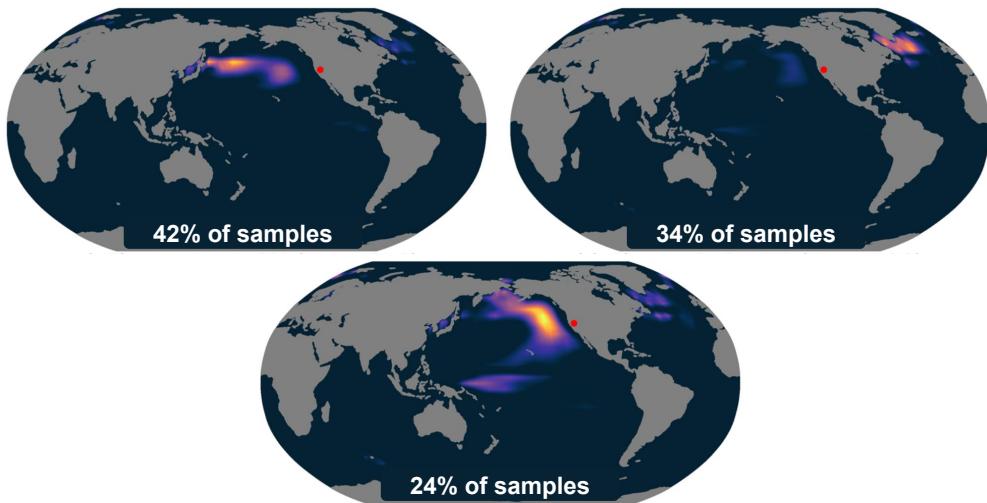
Predicting 5-year average surface temperature at each grid point
Applied to 1200 years of climate model CESM2 control simulation



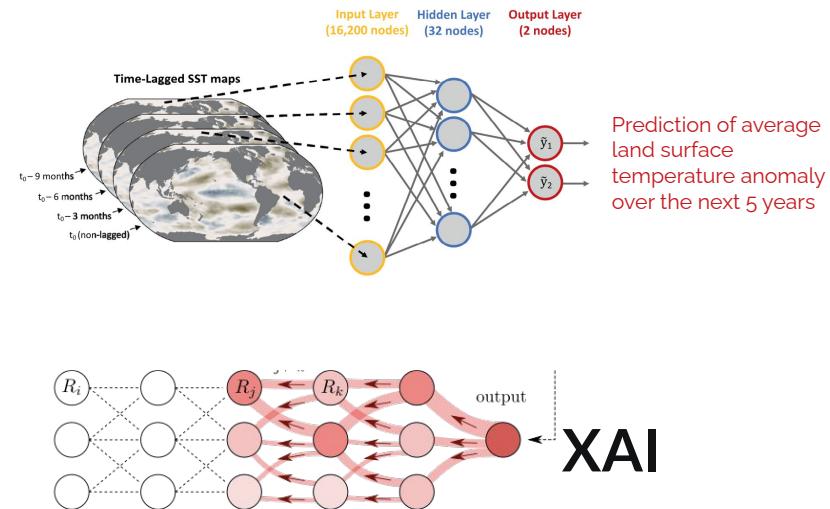
XAI reveals sources of predictability that vary in time and space



Toms, Barnes & Hurrell (2021)



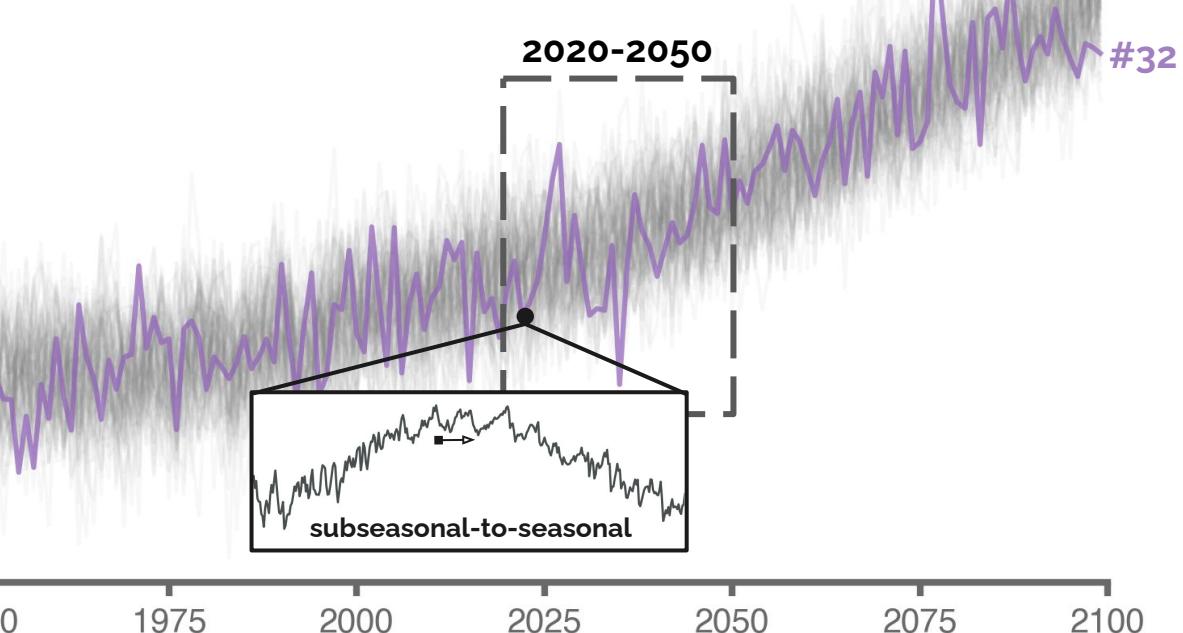
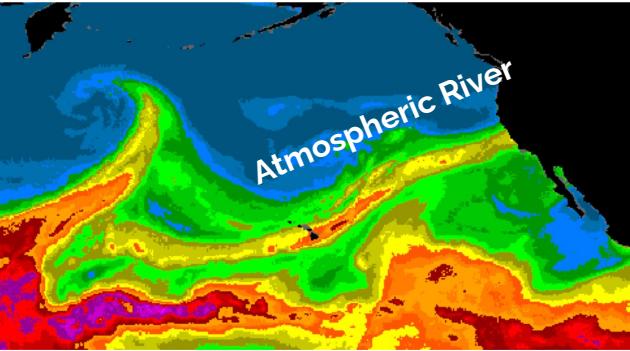
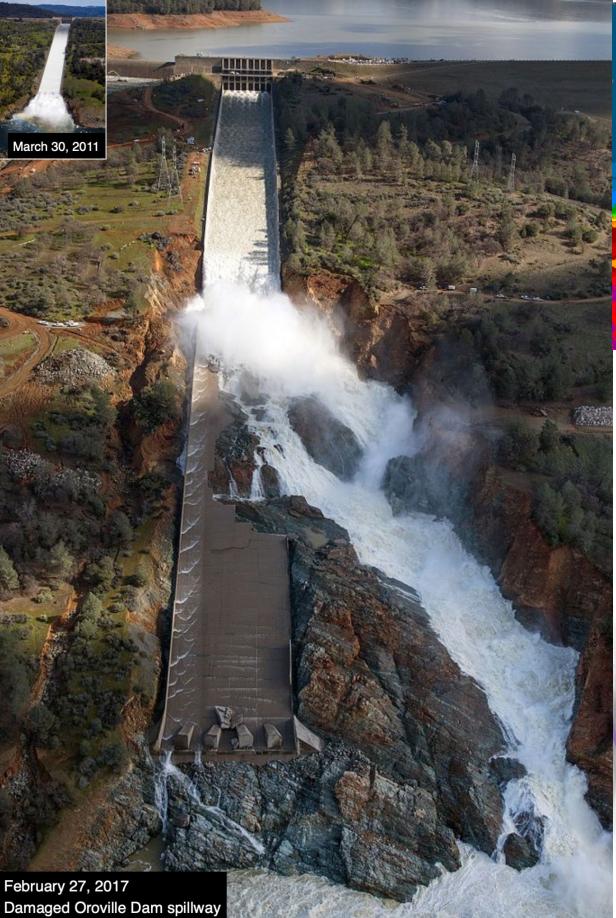
Predicting 5-year average surface temperature at each grid point
Applied to 1200 years of climate model CESM2 control simulation



XAI reveals sources of predictability that vary in time and space



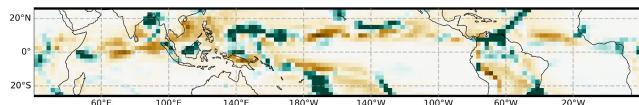
Toms, Barnes & Hurrell (2021)



Surface temperature over Chicago, IL
MPI-ESM Large Ensemble; historical + RCP8.5

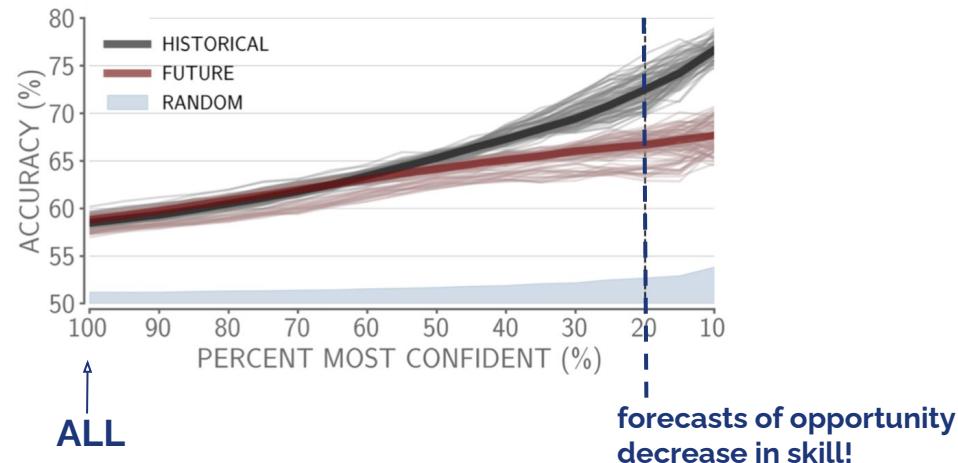
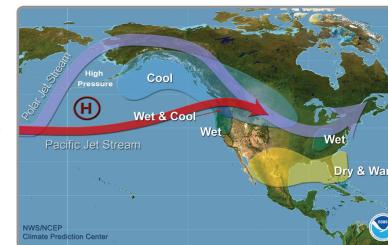
Input: Daily tropical precipitation

Trained on climate model **CESM2 [800 years of daily data]**



AI

Output: Pacific circulation 3 weeks later



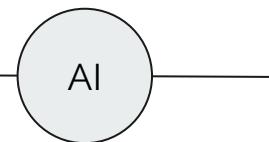
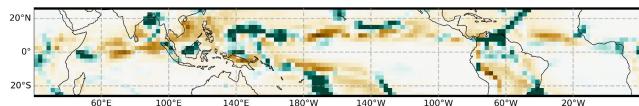
XAI allows us to quantify predictability in past and future climates and assess its sources.



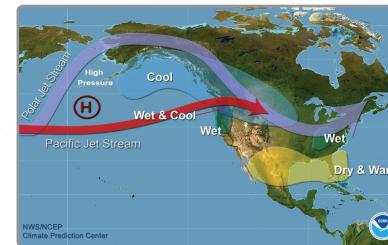
Mayer & Barnes (2021, 2022)

Input: Daily tropical precipitation

Trained on climate model **CESM2 [800 years of daily data]**

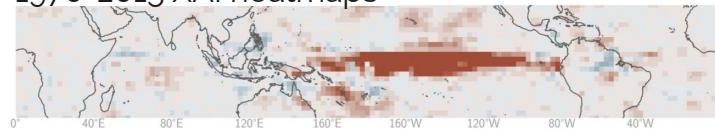


Output: Pacific circulation 3 weeks later

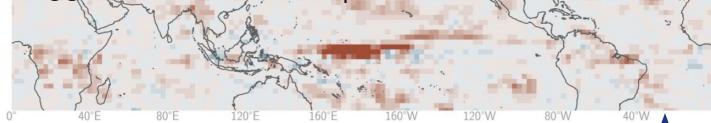


Most Relevant Regions

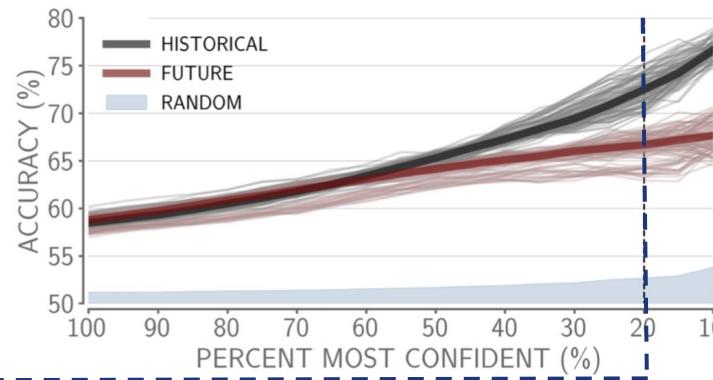
1970-2015 XAI heatmaps



2055-2100 XAI heatmaps



XAI



XAI allows us to quantify predictability in past and future climates and assess its sources.



Mayer & Barnes (2021, 2022)

The Earth system is coupled and complex.

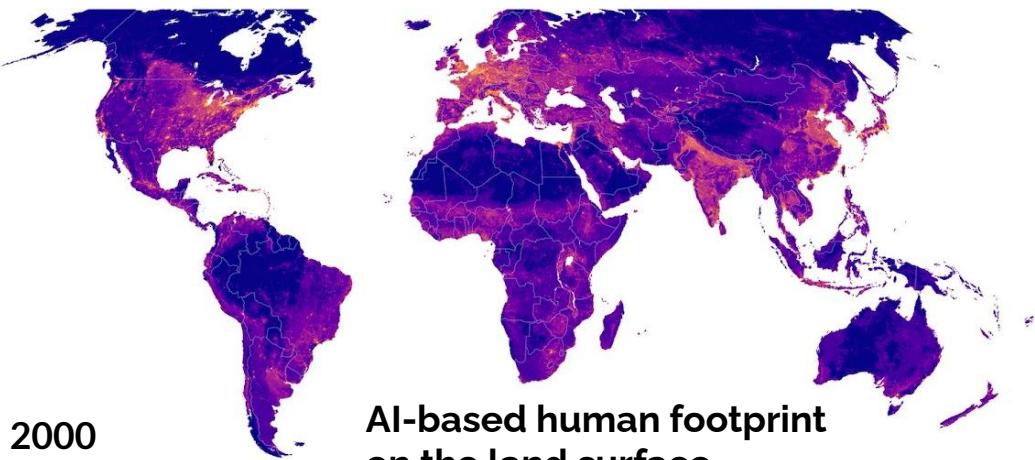
Human development and environmental factors cannot be easily decoupled.

Can we quantify the present and anticipate future changes in human development under future climate drivers?



Human-Earth System Interactions

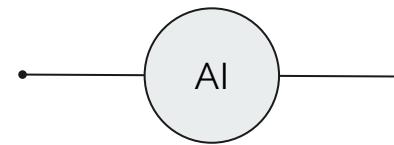
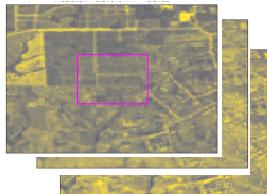
Quantify the present and predict the future transformations of the earth system by humanity in the presence of climate hazards



2000

AI-based human footprint
on the land surface

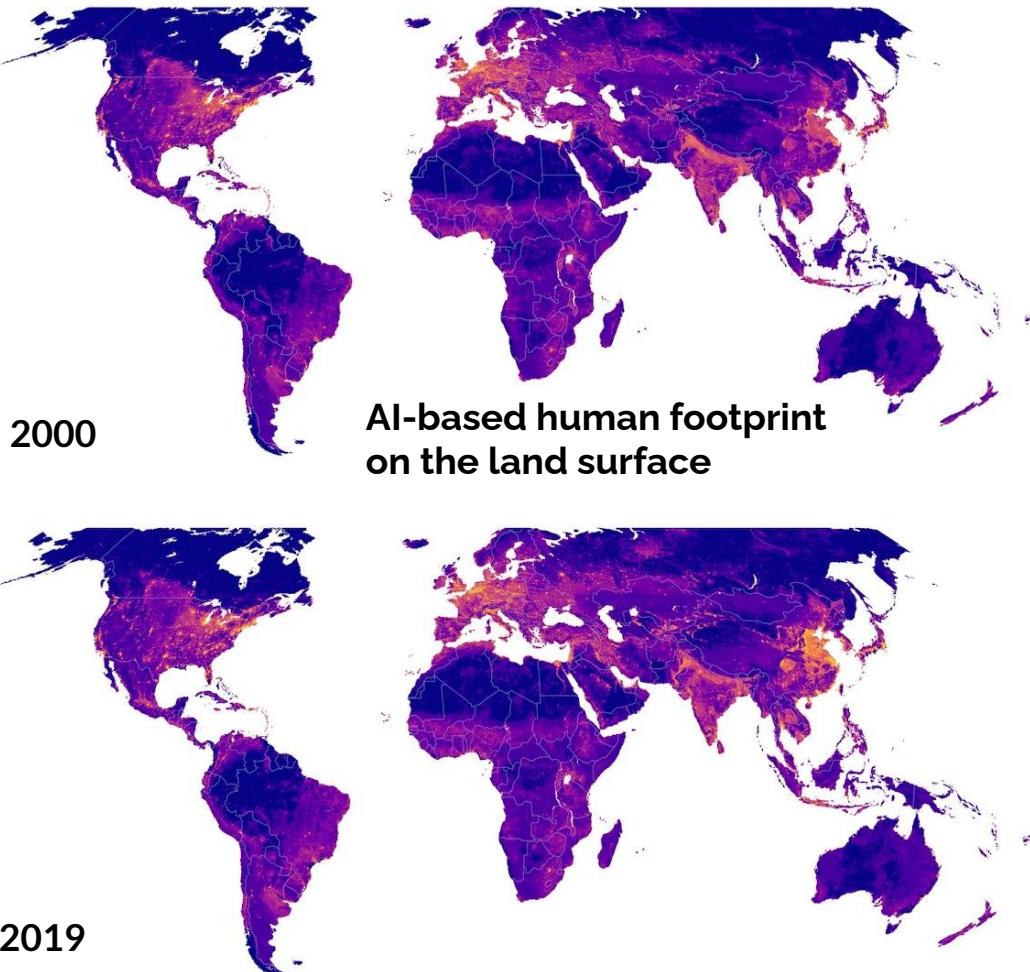
Landsat imagery from year 2000
3 bands x 120 pixels x 120 pixels



Keys, Barnes & Carter (2021)

Human-Earth System Interactions

Quantify the present and predict the future transformations of the earth system by humanity in the presence of climate hazards



Keys, Barnes & Carter (2021)

Human-Earth System Interactions

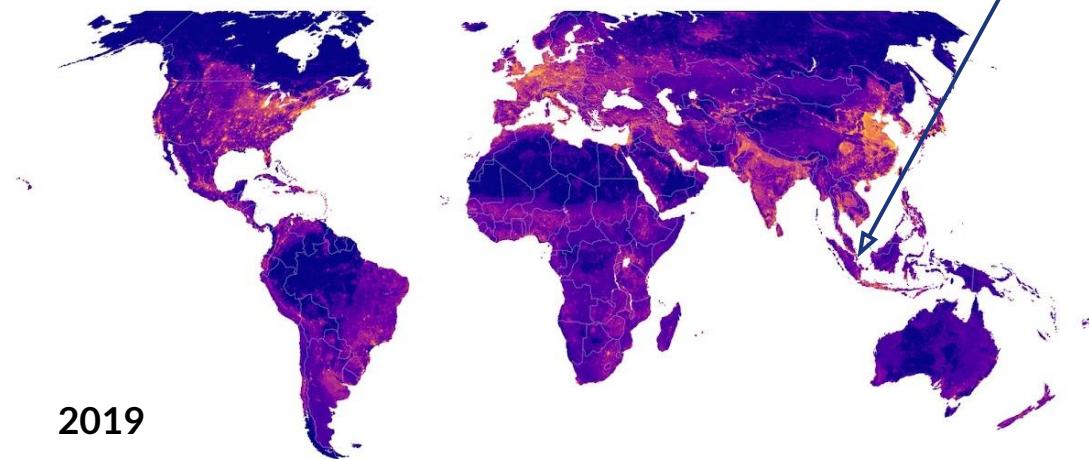
Quantify the present and predict the future transformations of the earth system by humanity in the presence of climate hazards



2000

2019

XAI



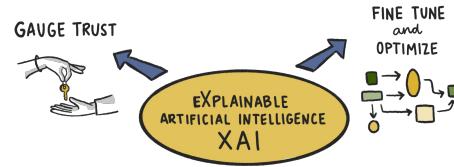
2019



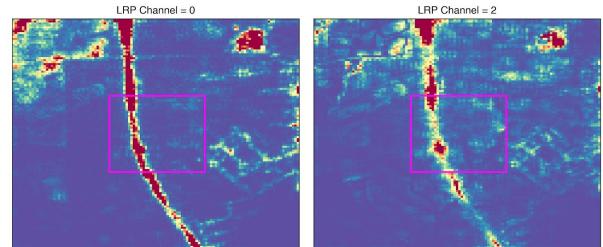
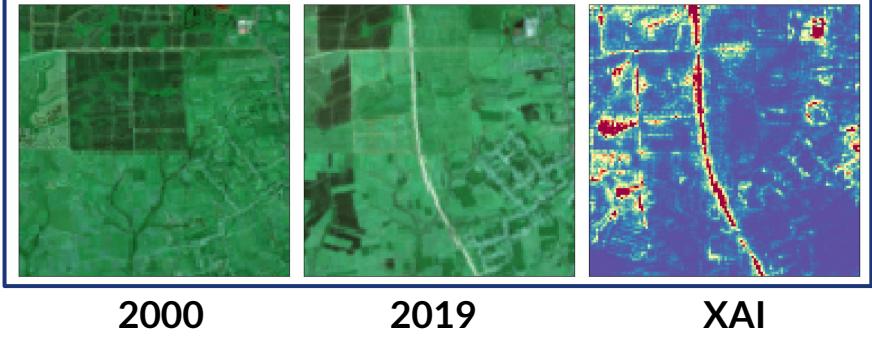
Keys, Barnes & Carter (2021)

Human-Earth System Interactions

Quantify the present and predict the future transformations of the earth system by humanity in the presence of climate hazards



XAI helped us determine which Landsat channels were most important for the network's decision; aiding us in refining our inputs.

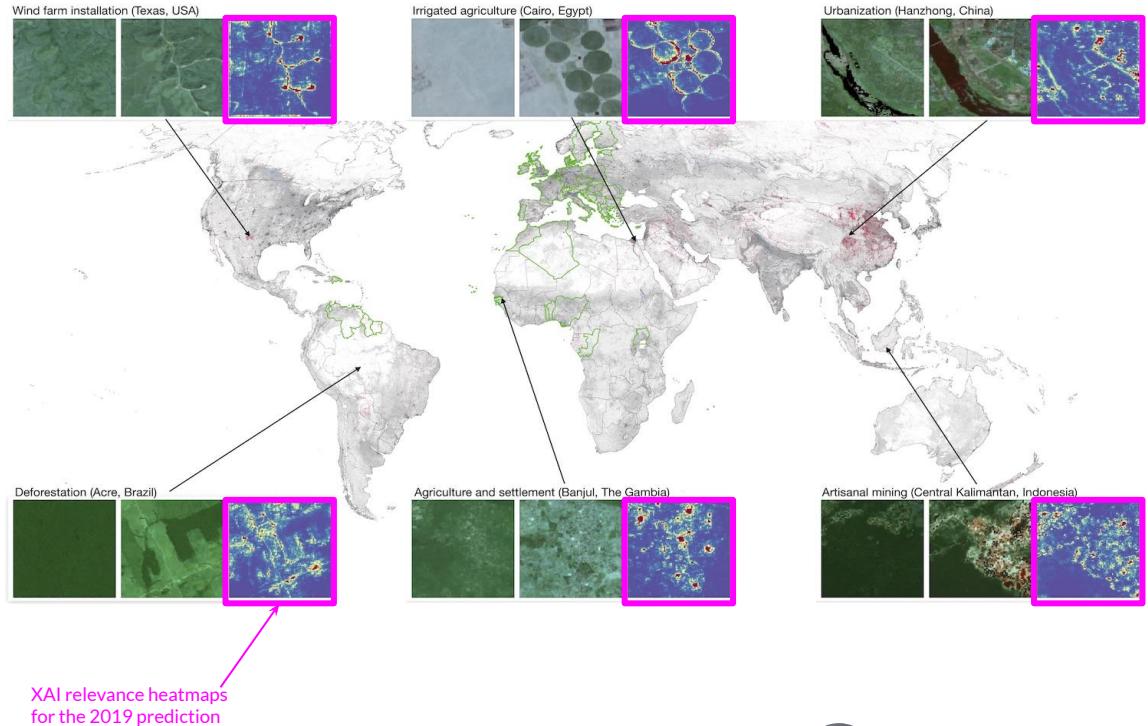


Keys, Barnes & Carter (2021)

Human-Earth System Interactions

Quantify the present and predict the future transformations of the earth system by humanity in the presence of climate hazards

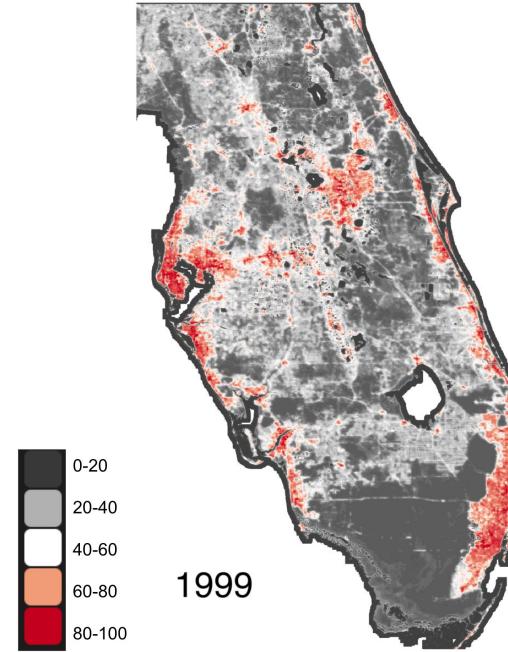
Changes in human footprint between 2000 and 2019



Keys, Barnes & Carter (2021)

Human-Earth System Interactions

Quantify the present and predict the future transformations of the earth system by humanity in the presence of climate hazards



Assess how climate hazards may interfere with global sustainable development in the coming years



Orihuela-Pinto, Keys, Davenport & Barnes (in prep)

Interpretable AI

AI models that are interpretable by design

- Explainable AI tells us *where*, but not *how*.
- Interpretable AI explicitly incorporates the decision-making process into its structure

PERSPECTIVE

<https://doi.org/10.1038/s42256-019-0048-x>

nature
machine intelligence

Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead

Cynthia Rudin 



Interpretable AI

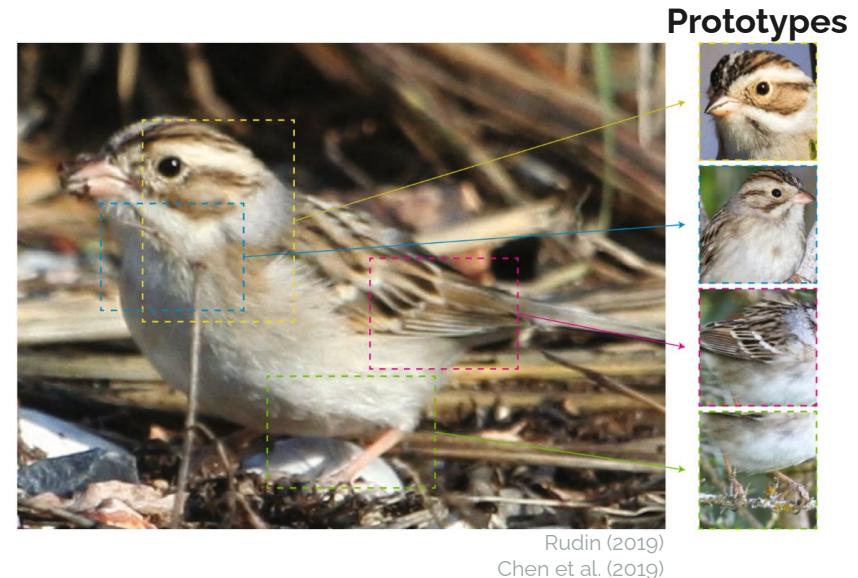
AI models that are interpretable by design

- Explainable AI tells us *where*, but not *how*.
- Interpretable AI **explicitly incorporates the decision-making process** into its structure

Our Current Goal:

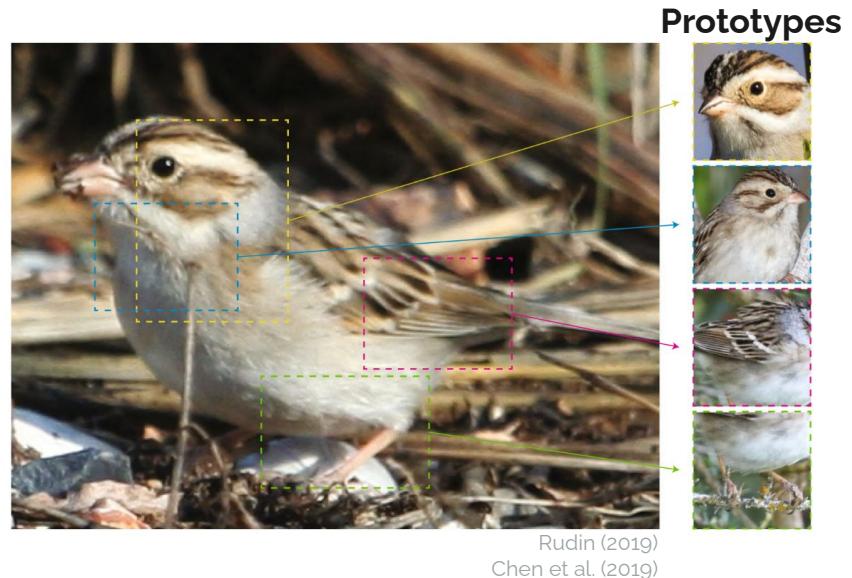
work toward building AI models that mimic scientific human reasoning to improve intrinsic interpretability

This Looks Like That



Building interpretable AI for climate applications

This Looks Like That Anywhere

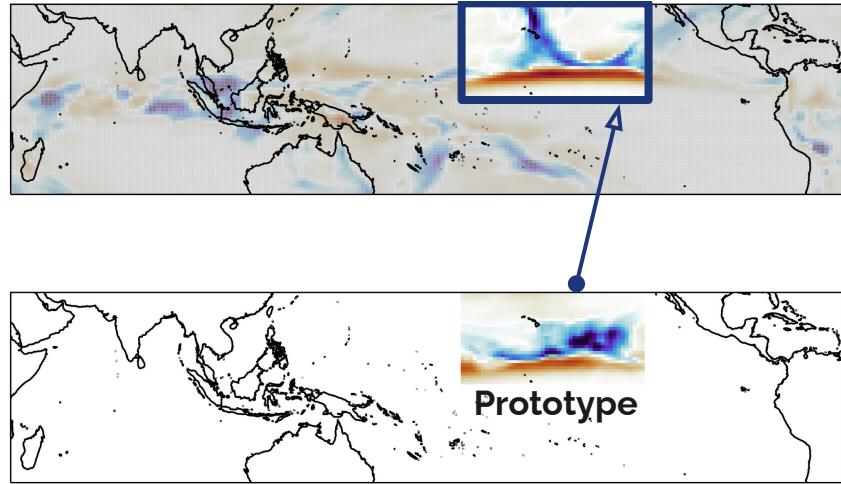


Building interpretable AI for climate applications

This Looks Like ***That*** ~~***Anywhere***~~
There

Building interpretable AI for climate applications

**This Looks Like ~~That Anywhere~~
There**



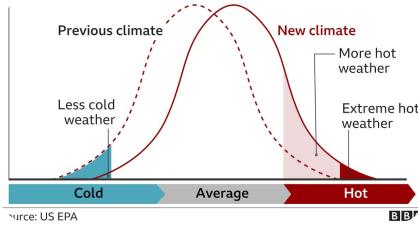
the network is constrained to make its predictions based on only #N prototypes – it needs to learn which ones

Building interpretable AI for climate applications

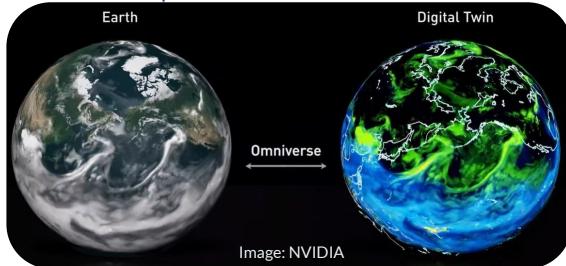
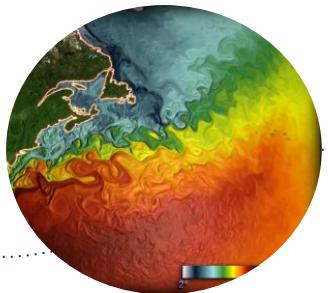


Gordillo and Barnes (under review)

Climate change influence the performance of AI models

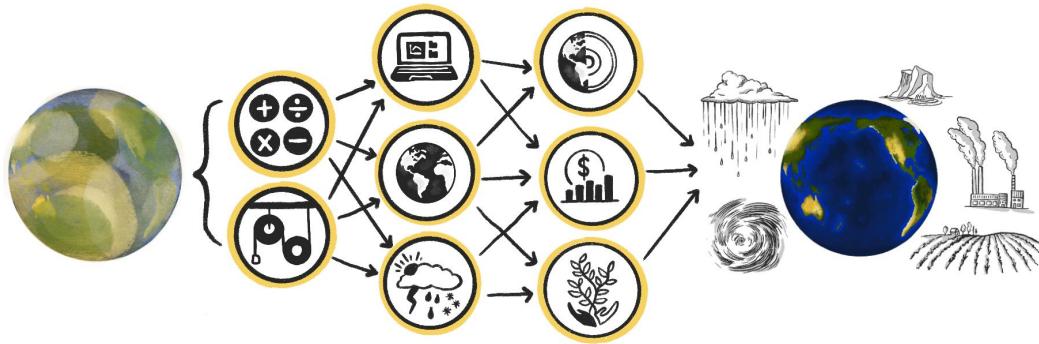


Role of AI in high-resolution Earth system modeling



AI for weather & climate model replacement

Additional New Frontiers for XAI + Climate



Thank you.

eabarnes@colostate.edu

<https://barnes.atmos.colostate.edu>

github: eabarnes1010