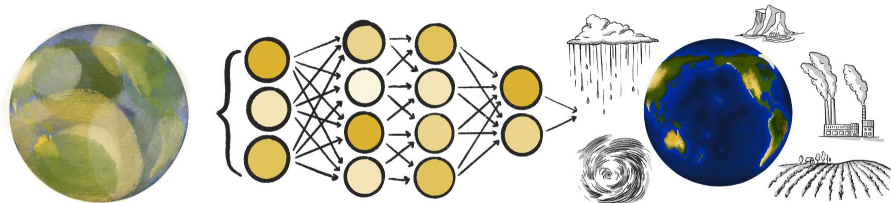# Explainable AI for Climate Science: Detection, Prediction and Discovery

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

email: eabarnes@colostate.edu
website: https://barnes.atmos.colostate.edu
twitter: @atmosbarnes
github: eabarnes1010

ATMOSPHERIC SCIENCE
COLORADO STATE UNIVERSITY

WCRP hybrid symposium on Frontiers in Subseasonal to Decadal Prediction
March 28, 2023

One of our jobs as scientists is to sift through piles of data and try to extract useful relationships that apply elsewhere, i.e. that are applicable "out of sample".

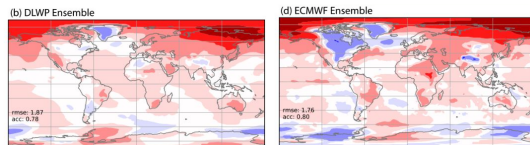This is what many machine learning methods are designed to do.

_____

*throughout this talk I will use "AI" and "machine learning" or "ML" interchangeably

# ML for S2D Predictability & Prediction

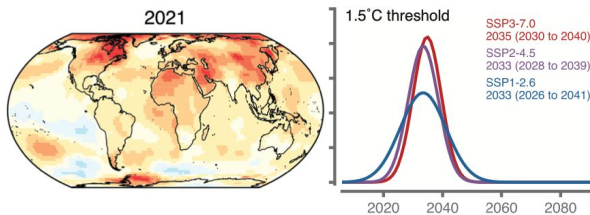There are many ways that ML can be applied to try and improve understanding of S2D predictability

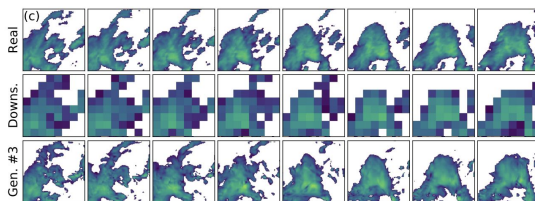## Deep Learning WP Models
e.g. *Weyn et al. (2021)*



## Post-processing of Multi-Model Ensembles
e.g. Haupt et al. (2021), Schumacher et al. (2021), Gronquist et al. (2020), *Diffenbaugh and Barnes (in review)*
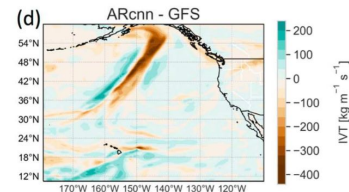


## Downscaling for Regional Impacts
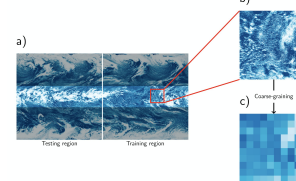e.g. *Leinonen et al. (2020)*



## Predicting the Errors of Forecast Systems
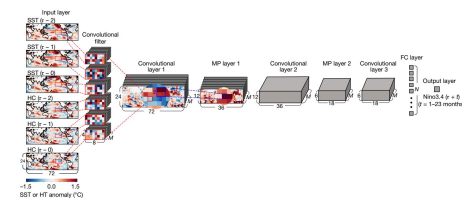e.g. *Chapman et al. (2019)*, Cahill et al. (in prep), Pan et al. (2021)



## 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); *Brenowitz and Bretherton (2018)*



## Statistical Model Predictions
e.g. Mayer and Barnes (2021); Hassanibesheli et al. (2022); *Ham et al. (2019)*

# Opening the Black Box with XAI

**More and more papers are coming out demonstrating the use of ML explainability methods for geoscience**



MAKING THE BLACK BOX MORE TRANSPARENT
Understanding the Physical Implications of Machine Learning

Amy McGovern, Ryan Lagerquist, David John Gagne II, G. Eli Jergensen, Kimberly L. Elmore, Cameron R. Homeyer, and Travis Smith

Machine learning model interpretation and visualization focusing on meteorological domains are introduced and analyzed.

JAMES | Journal of Advances in Modeling Earth Systems

RESEARCH ARTICLE

Physically Interpretable Neural Networks for the Geosciences: Applications to Earth System Variability

Benjamin A. Toms, Elizabeth A. Barnes, and Imme Ebert-Uphoff

arXiv > physics > arXiv:2103.10005

Physics > Geophysics

[Submitted on 18 Mar 2021]

**Neural Network Attribution Methods for Problems in Geoscience: A Novel Synthetic Benchmark Dataset**

Antonios Mamalakis, Imme Ebert-Uphoff, Elizabeth A. Barnes

arXiv > physics > arXiv:2202.03407

Physics > Geophysics

[Submitted on 7 Feb 2022]

**Investigating the fidelity of explainable artificial intelligence methods for applications of convolutional neural networks in geoscience**

Antonios Mamalakis, Elizabeth A. Barnes, Imme Ebert-Uphoff

BAMS | ISSUES | EARLY ONLINE RELEASE | COLLECTIONS | FOR AUTHORS

RESEARCH ARTICLE | 31 AUGUST 2020

**Evaluation, Tuning and Interpretation of Neural Networks for Working with Images in Meteorological Applications**

Imme Ebert-Uphoff | Kyle Hilburn
Bull. Amer. Meteor. Soc. 1–49.

# XAI Attribution Methods

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



e.g. Montavon et al. (2017), Pattern Recognition; Montavon et al. (2018), Digital Signal Processing

# 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.

GAUGE TRUST

FINE TUNE and OPTIMIZE

EXPLAINABLE ARTIFICIAL INTELLIGENCE XAI

LEARN NEW SCIENCE

*Our ultimate goal is to understand causality, and XAI does not give us this. But it is a step in the right direction.*

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

*e.g. Labe and Barnes (2022), Diffenbaugh and Barnes (2023), Rader et al. (2022), Labe and Barnes (2021), Barnes et al. (2020a), Barnes et al. (2019)*

global warming projections

1900                    2000                    2100

Surface temperature over Fort Collins, CO
CanESM2 simulation
historical + SSP3-7.0

# Time Remaining Until Critical Warming Thresholds are Reached



2011–2020 was
around 1.1°C warmer
than 1850–1900

the last time global surface temperature was sustained
at or above 2.5°C was over 3 million years ago

The world at
+1.5°C

The world at
+2°C

The world at
+3°C

The world at
+4°C

0

1

Global warming level (GWL) above 1850-1900

°C

? ? ? ?

# Time Remaining Until Critical Warming Thresholds are Reached

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

years until warming threshold is reached + uncertainties

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

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

years until warming threshold is reached + uncertainties

Climate Model Results

- training
- validation
- testing

testing MAE = 2.7 yrs.

true number of years

**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)

1980    1992

2016    2021

**Observations**
Berkeley Earth Surface Temperature

−4   −2   0   2   4

temperature anomaly (relative to 1951-1980)

# 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

1980     1992     2016     2021

temperature anomaly (relative to 1951-1980)

2.0°C threshold

Pinatubo eruption

years until treshold

SSP3-7.0

slope = -1.0 years/yr

year

**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)

1980     1992     2016     2021

temperature anomaly (relative to 1951-1980)

higher likelihood of reaching 2°C in the Low scenario than indicated in some previous assessments—although the possibility it could be avoided is not ruled out.

2.0°C threshold

SSP3-7.0
2050 (2043 to 2058)

SSP2-4.5
2049 (2043 to 2055)

SSP1-2.6
2054 (2044 to 2065)

year

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

A. CMIP6 10-15 yrs to threshold

B. Observations 2018-2021

XAI

further from threshold — nearer to threshold

**Use the trained AI model to predict thresholds based on maps of the observed climate**
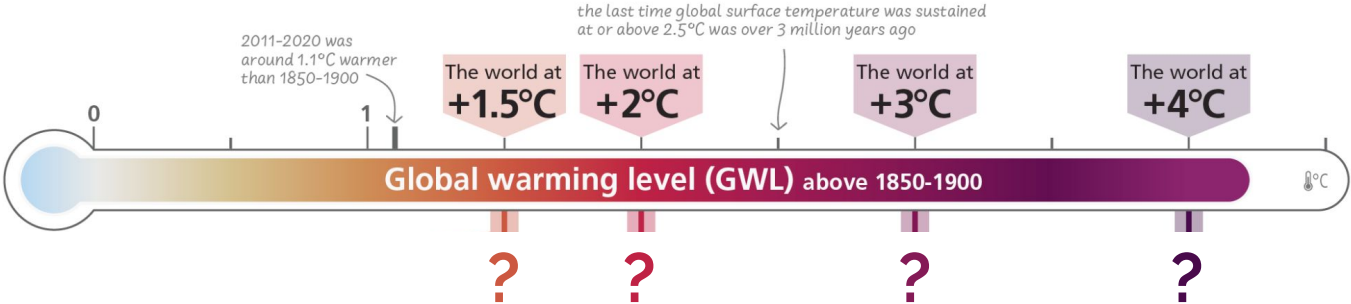
# **AI to leverage** imperfect climate models to better **constrain future projections** by **fusing** simulations and observations.

*e.g. Labe and Barnes (2022), Diffenbaugh and Barnes (2023), Rader et al. (2022), Labe and Barnes (2021), Barnes et al. (2020a), Barnes et al. (2019)*
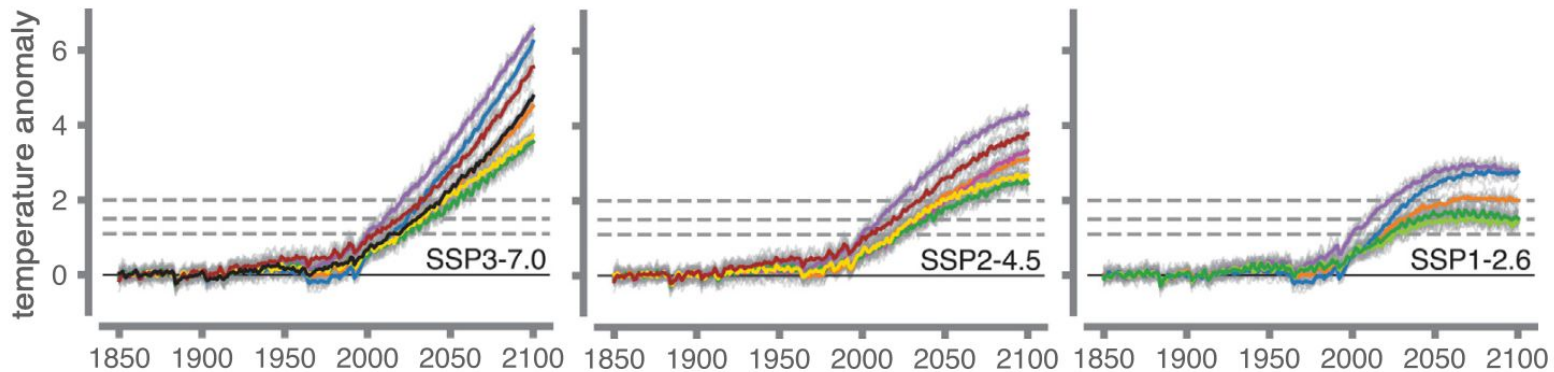
global warming projections

1900    2000    2100

Surface temperature over Fort Collins, CO
CanESM2 simulation
historical + SSP3-7.0

# AI to explore earth system predictability weeks-to-years in advance to study prediction, dynamics, and change.

*e.g. Gordon and Barnes (2022), Labe and Barnes (2022), Gordon, Barnes and Hurrell (2021), Toms, Barnes and Hurrell (2021), Mayer and Barnes (2022), Mayer and Barnes (2021), Barnes et al. (2020), Barnes et al. (2020)*

multi-year predictability

1900          2000          2100

Surface temperature over Fort Collins, CO
CanESM2 simulation
historical + SSP3-7.0

**AI to explore** earth system **predictability weeks-to-years in advance** to study prediction, dynamics, and change.

*e.g. Gordon and Barnes (2022), Labe and Barnes (2022), Gordon, Barnes and Hurrell (2021), Toms, Barnes and Hurrell (2021), Mayer and Barnes (2022), Mayer and Barnes (2021), Barnes et al. (2020), Barnes et al. (2020)*

subseasonal-to-seasonal

1900          2000          2100

Surface temperature over Fort Collins, CO
CanESM2 simulation
historical + SSP3-7.0

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.

80% of the data is noise

20% of the data is signal

**Impossible predictions may be hampering AI learning of predictable behaviour**

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.**



Madden-Julian oscillation
[~30 days]

El Nino Southern Oscillation
[~2 years]

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

prediction of
future regional
sea surface
temperature

**thoughtful choices of
inputs and outputs can
allow attribution of
sources of predictability**

# Attributing external + internal sources of predictability

**External Forcing Model**

**Internal Variability Model**

**Full Model Accuracy**

accuracy (%)

**Attributing external + internal sources of predictability**

Gordon & Barnes (in prep)

**External Forcing Model**

**Internal Variability Model**

**Full Model Accuracy**

accuracy (%)

**Full Model − External Only Model =**

# Attributing external + internal sources of predictability

Gordon & Barnes (in prep)

**External Forcing Model**

**Internal Variability Model**

Full Model Accuracy

Skill Added by Internal Variability

**Attributing external + internal sources of predictability**

Gordon & Barnes (in prep)

**External Forcing Model**

**Internal Variability Model**
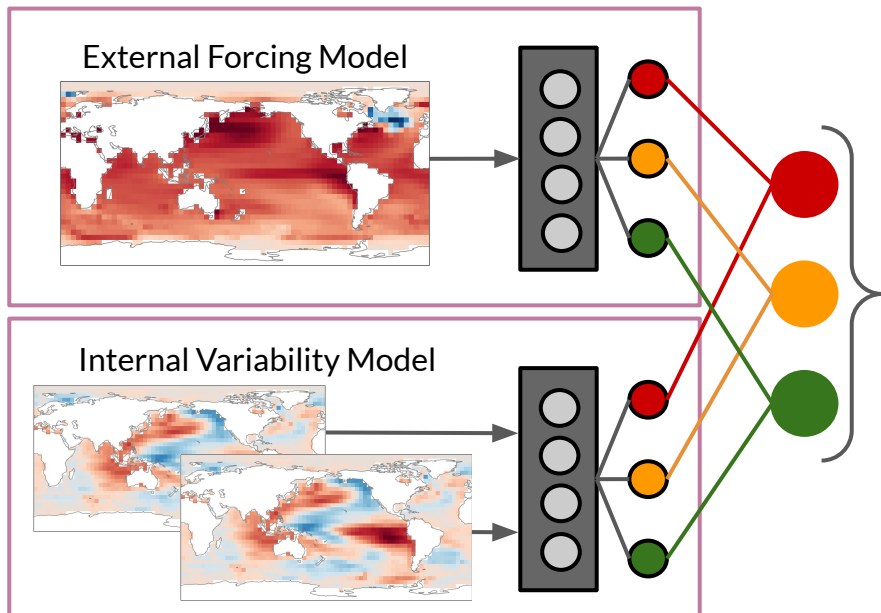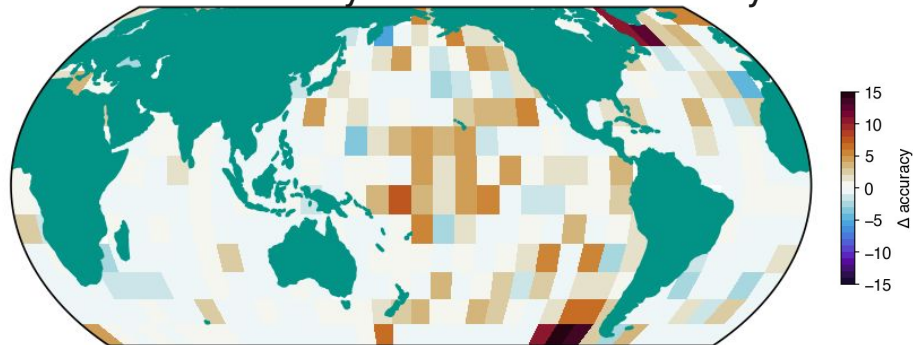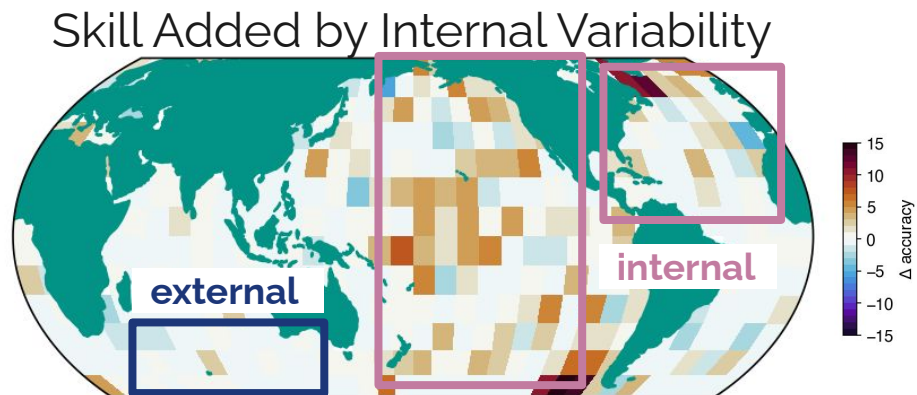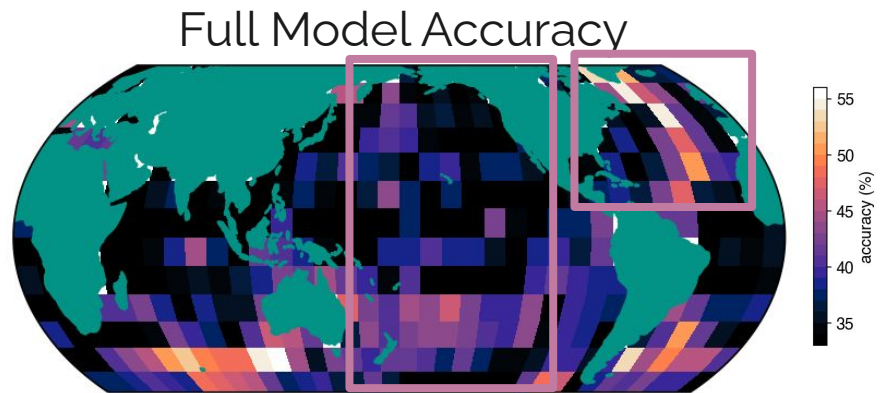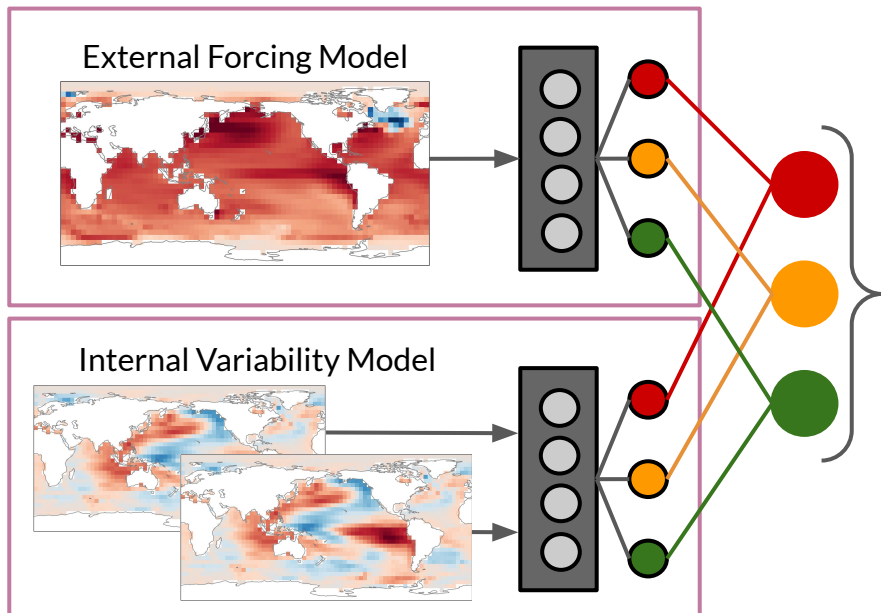
Full Model Accuracy

Skill Added by Internal Variability

external

internal

**Attributing external + internal sources of predictability**

Gordon & Barnes (in prep)

**past sea-surface temperatures**

3-8 years before

2-3 years before

1-2 years before

0-1 years before

**future sea surface temperature anomalies***
*for one point*
[0-5 years]

**positive**

neutral

**negative**

*can predict a range of variables

# Predict ocean temperatures 5 years later

Davenport & Barnes (in prep)

# CLIMATE MODEL DATA

**Overall Accuracy**



0.4    0.6    0.8    1.0

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 (in prep)

# CLIMATE MODEL DATA

### Overall Accuracy

### Accuracy for 40% most confident predictions



0.4    0.6    0.8    1.0

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

# CLIMATE MODEL DATA

**Overall Accuracy**

**Accuracy for 40% most confident predictions**

**Accuracy for 20% most confident predictions**
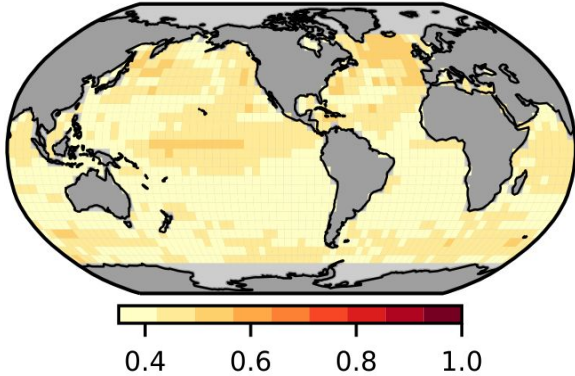


0.4   0.6   0.8   1.0

>60-80% accuracy!

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

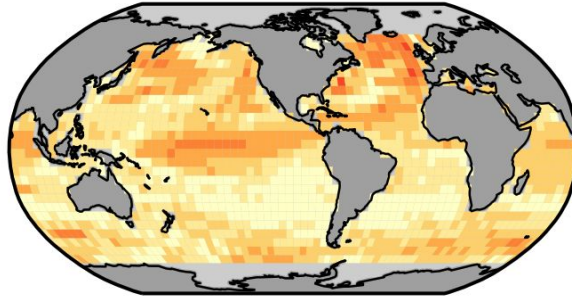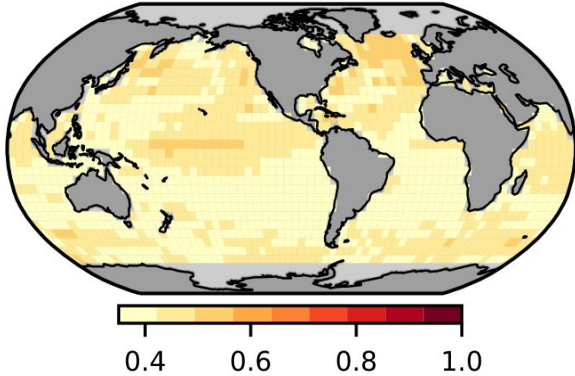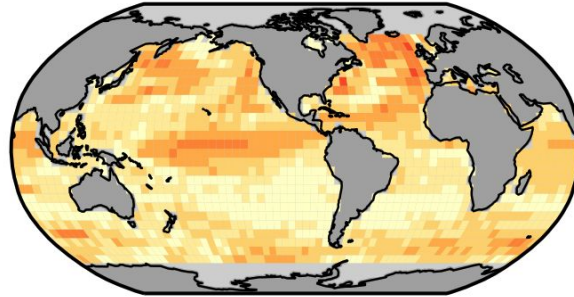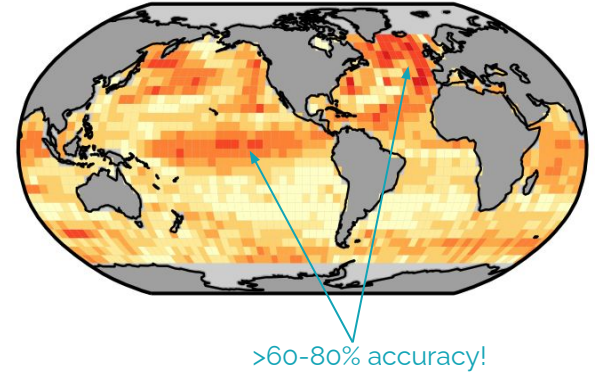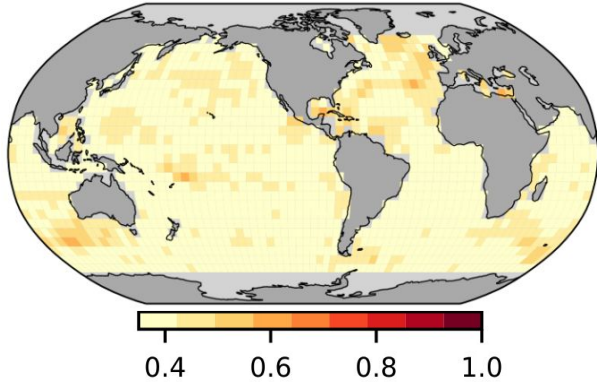## Focusing on when the AI is most confident leads to accurate predictions

Davenport & Barnes (in prep)

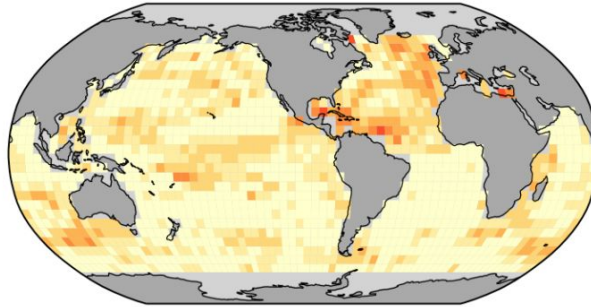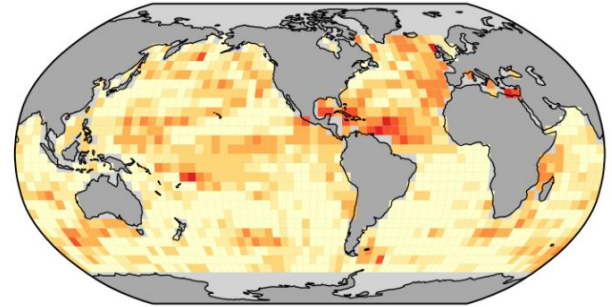# OBSERVATIONS

**Overall Accuracy**

**Accuracy for 40% most confident predictions**

**Accuracy for 20% most confident predictions**
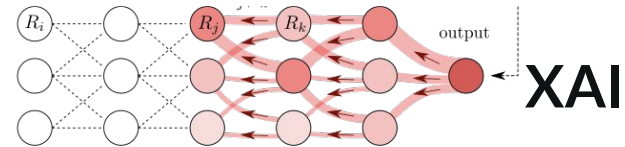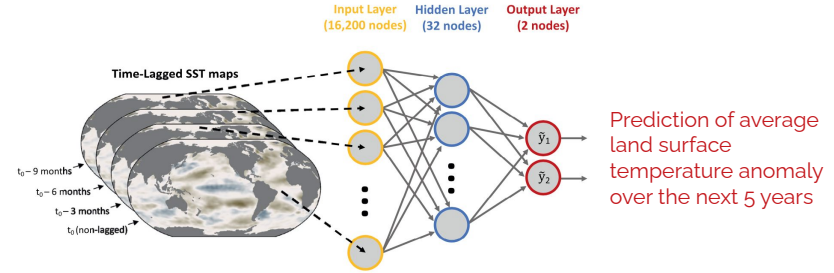


0.4   0.6   0.8   1.0

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 (in prep)

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

Input Layer (16,200 nodes)
Hidden Layer (32 nodes)
Output Layer (2 nodes)

Time-Lagged SST maps

$t_0 - 9$ months
$t_0 - 6$ months
$t_0 - 3$ months
$t_0$ (non-lagged)

$\hat{y}_1$
$\hat{y}_2$

Prediction of average land surface temperature anomaly over the next 5 years

$R_i$
$R_j$
$R_k$
output

XAI

**XAI reveals sources of predictability that vary in time and space**
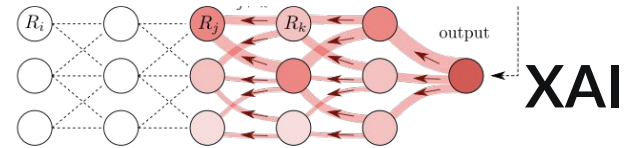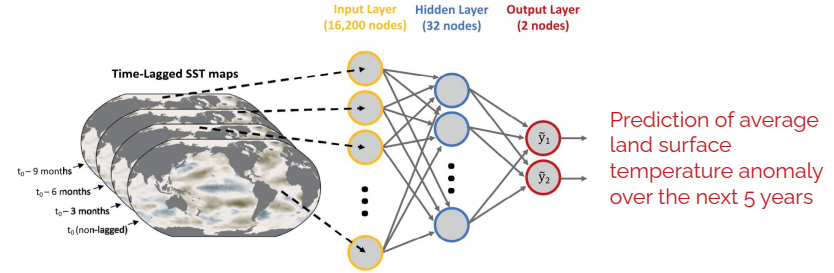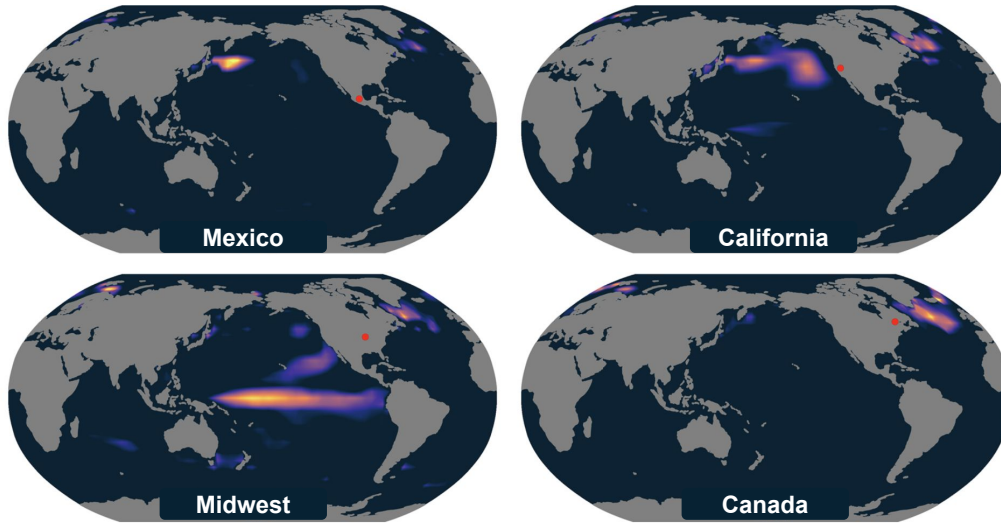
Predicting 5-year average surface temperature at each grid point
Applied to 1200 years of CESM2 control simulation

Input Layer (16,200 nodes)   Hidden Layer (32 nodes)   Output Layer (2 nodes)

Time-Lagged SST maps

$t_0 - 9$ months
$t_0 - 6$ months
$t_0 - 3$ months
$t_0$ (non-lagged)

$\hat{y}_1$
$\hat{y}_2$

Prediction of average land surface temperature anomaly over the next 5 years

Mexico

California

Midwest

Canada

$R_i$   $R_j$   $R_k$   output
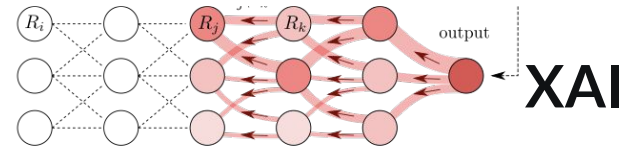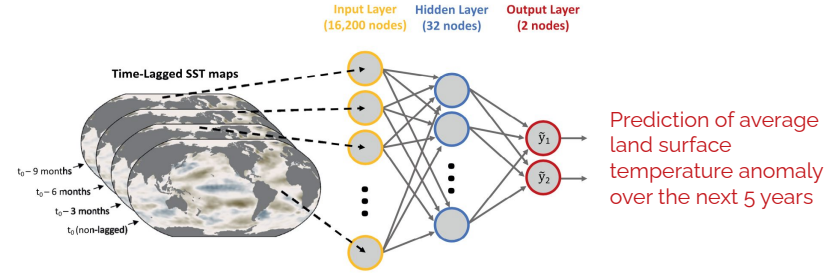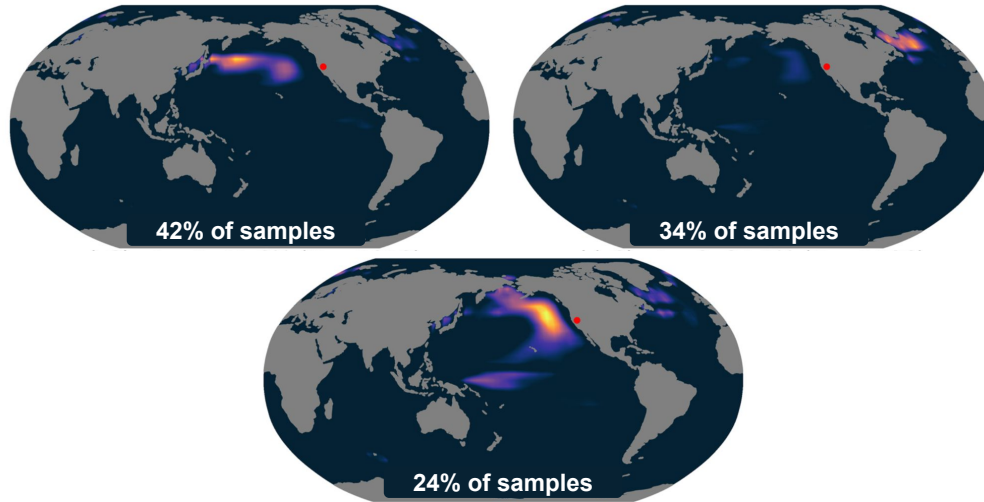
XAI

**XAI reveals sources of predictability that vary in time and space**

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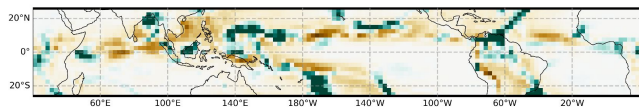Applied to 1200 years of CESM2 control simulation

Input Layer
(16,200 nodes)
Hidden Layer
(32 nodes)
Output Layer
(2 nodes)

Time-Lagged SST maps

$t_0 - 9$ months
$t_0 - 6$ months
$t_0 - 3$ months
$t_0$ (non-lagged)

$\hat{y}_1$
$\hat{y}_2$

Prediction of average
land surface
temperature anomaly
over the next 5 years

42% of samples

34% of samples

24% of samples

$R_i$    $R_j$    $R_k$    output

XAI

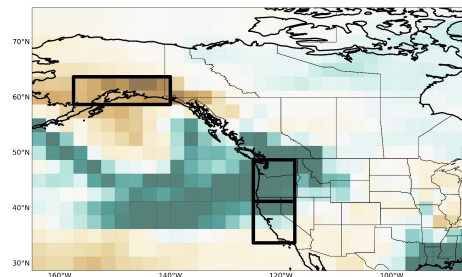**XAI reveals sources of predictability that vary in time and space**

**Input:** Daily tropical precipitation

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

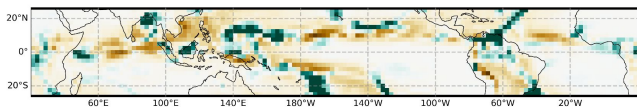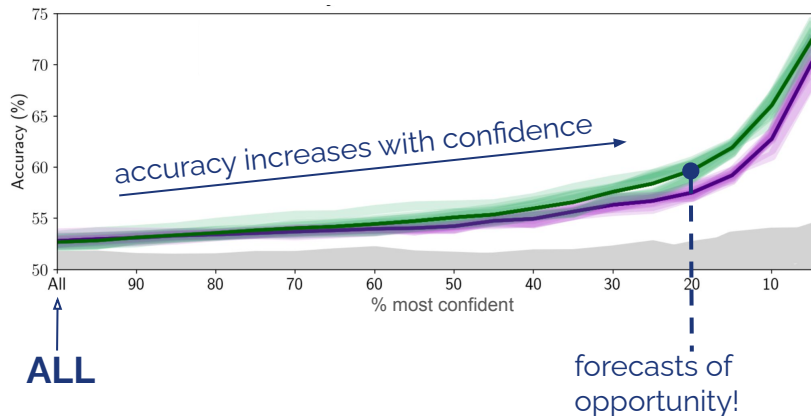**Output:** Precipitation 3-4 weeks later

AI

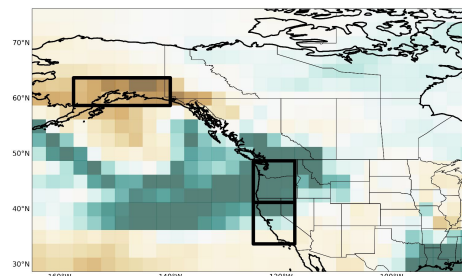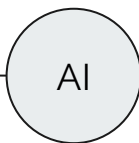**XAI allows us to quantify predictability in past and future climates and assess changes in sources.**

**Input:** Daily tropical precipitation

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

**Output:** Precipitation 3-4 weeks later



AI

accuracy increases with confidence
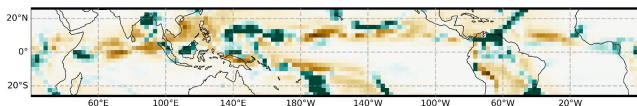
**ALL**

forecasts of opportunity!

**XAI allows us to quantify predictability in past and future climates and assess changes in sources.**
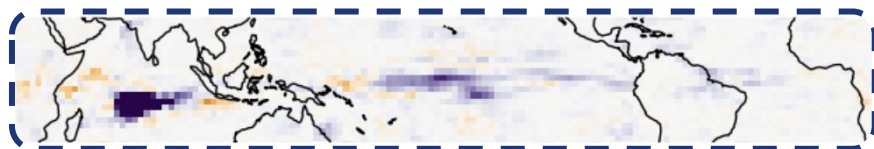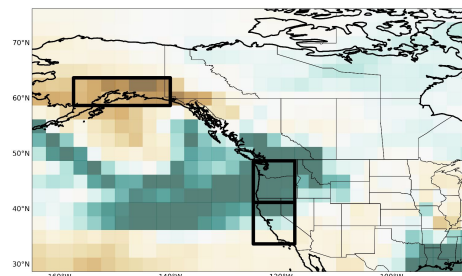
**Input:** Daily tropical precipitation

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

**Output:** Precipitation 3-4 weeks later



AI



**XAI**

XAI allows us to quantify predictability in past and future climates and assess changes in sources.
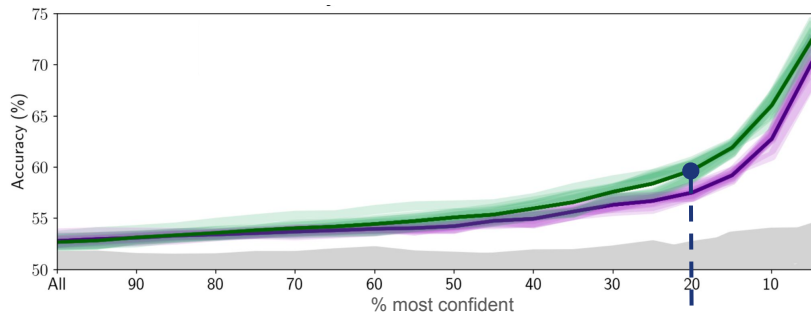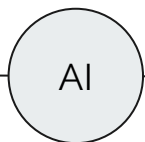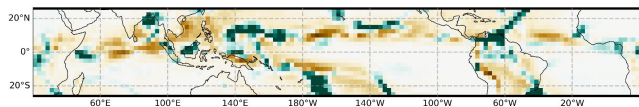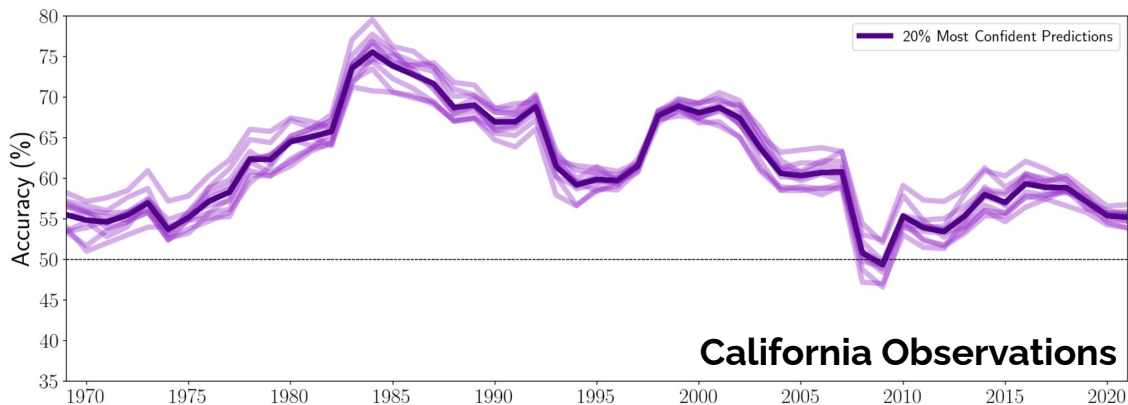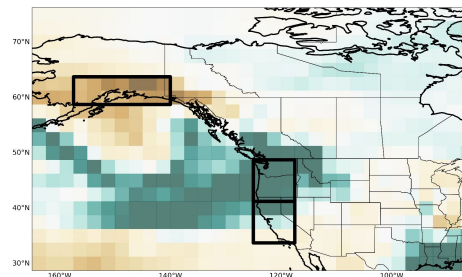
**Input:** Daily tropical precipitation
Trained on climate model **CESM2 [800 years of daily data]**

**Output:** Precipitation 3-4 weeks later

AI

California Observations
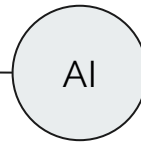
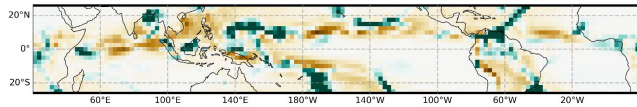20% Most Confident Predictions

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

**Input:** Daily tropical precipitation

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

**Output:** Pacific circulation 3 weeks later



AI

ALL

forecasts of opportunity
decrease in skill!

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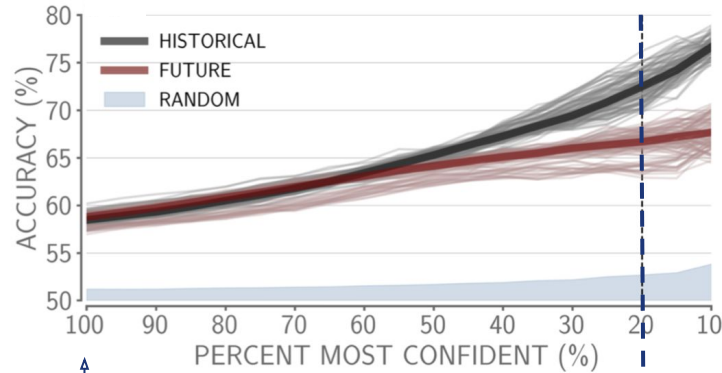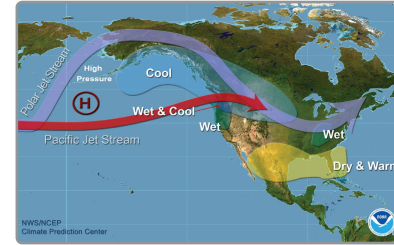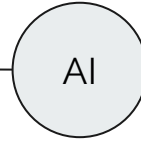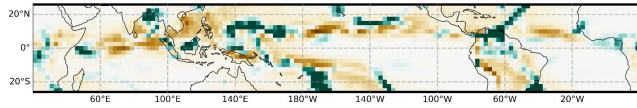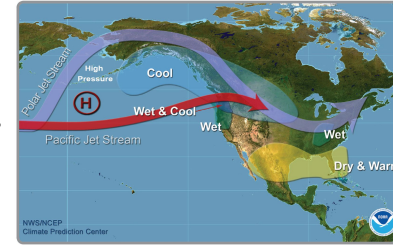and assess changes in sources.**

Mayer & Barnes (2021, 2022)

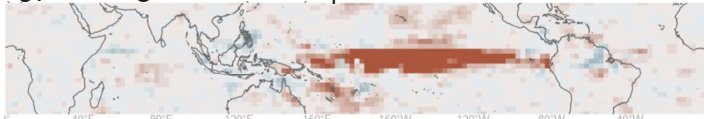**Input:** Daily tropical precipitation

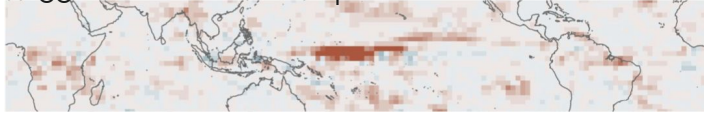Trained on climate model **CESM2 [800 years of daily data]**

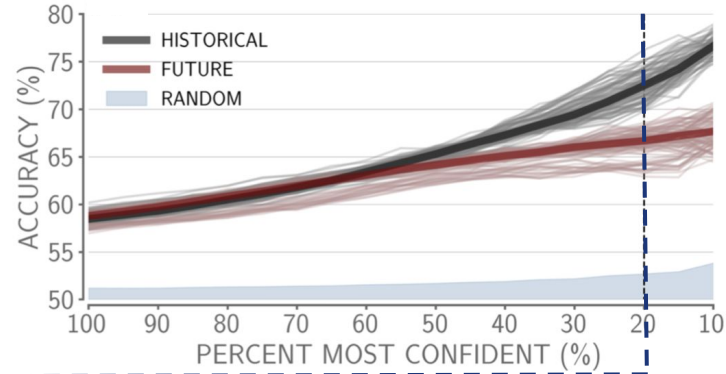**Output:** Pacific circulation 3 weeks later



1970-2015 XAI heatmaps

2055-2100 XAI heatmaps

XAI

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

Mayer & Barnes (2021, 2022)

The future of actionable climate predictions requires the **mixing of knowledge**.

And ultimately we want more than just a prediction - we want to know "**why?**"

Explainable AI has a lot to offer climate prediction.

The future of actionable climate predictions requires the **mixing of knowledge**.

And ultimately we want more than just a prediction - we want to know "**why?**"

# Thank you.

eabarnes@colostate.edu

Explainable AI has a lot to offer climate prediction.

The future of actionable climate predictions requires the **mixing of knowledge**.

And ultimately we want more than just a prediction - we want to know "**why?**"