

Geophysical Research Letters[®]



RESEARCH LETTER

10.1029/2023GL106278

Key Points:

- Neural networks outperform persistence forecasts in predicting extreme states of North Atlantic sea surface temperature out to 25 years
- An explainable neural network technique reveals successful predictions rely consistently on the Transition Zone region
- Neural networks trained on climate model output predict the phasing of multidecadal variability on an observation-based data set

Supporting Information:

Supporting Information may be found in the online version of this article.

Correspondence to:




G. Liu,
glennliu@mit.edu

Citation:

Liu, G., Wang, P., & Kwon, Y.-O. (2023). Physical insights from the multidecadal prediction of North Atlantic sea surface temperature variability using explainable neural networks. *Geophysical Research Letters*, 50, e2023GL106278. <https://doi.org/10.1029/2023GL106278>

Received 13 SEP 2023
Accepted 26 NOV 2023

Physical Insights From the Multidecadal Prediction of North Atlantic Sea Surface Temperature Variability Using Explainable Neural Networks

Glenn Liu^{1,2} , Peidong Wang³ , and Young-Oh Kwon² 

¹MIT-WHOI Joint Program in Oceanography/Applied Ocean Science and Engineering, Cambridge, MA, USA, ²Physical Oceanography Department, Woods Hole Oceanographic Institution, Falmouth, MA, USA, ³Department of Earth, Atmospheric, and Planetary Sciences, Massachusetts Institute of Technology, Cambridge, MA, USA

Abstract North Atlantic sea surface temperatures (NASST), particularly in the subpolar region, are among the most predictable in the world's oceans. However, the relative importance of atmospheric and oceanic controls on their variability at multidecadal timescales remain uncertain. Neural networks (NNs) are trained to examine the relative importance of oceanic and atmospheric predictors in predicting the NASST state in the Community Earth System Model 1 (CESM1). In the presence of external forcings, oceanic predictors outperform atmospheric predictors, persistence, and random chance baselines out to 25-year leadtimes. Layer-wise relevance propagation is used to unveil the sources of predictability, and reveal that NNs consistently rely upon the Gulf Stream-North Atlantic Current region for accurate predictions. Additionally, CESM1-trained NNs successfully predict the phasing of multidecadal variability in an observational data set, suggesting consistency in physical processes driving NASST variability between CESM1 and observations.

Plain Language Summary North Atlantic sea surface temperatures, particularly in the subpolar region, are among the most predictable locations in the world's oceans. However, it remains uncertain if processes in the atmosphere or ocean are more important for driving temperature fluctuations in this region occurring over multiple decades. We use a machine learning approach to predict the sea surface temperature state from climate model outputs, given snapshots of atmospheric or oceanic variables. Ocean variables lead to more accurate predictions relative to atmospheric variables and standard prediction baselines out to 25 years ahead if processes that drive the trends in climate, such as human-induced warming, are present in the data. These successful predictions arise consistently from the same region near the Gulf Stream-North Atlantic Current region. Despite being trained on climate models, the neural networks can predict the timing of observed positive and negative states of real-world sea surface temperatures, suggesting that there is potential for using model output to train neural networks at predicting the actual North Atlantic sea surface variability.

1. Introduction

Sea surface temperature (SST) anomalies averaged over the North Atlantic region exhibit alternating warm and cold periods on multidecadal timescales, known as the Atlantic Multidecadal Variability (AMV) or Atlantic Multidecadal Oscillation. The societal relevance of predicting AMV is underscored by linkages to multidecadal variations across multiple Earth system processes both within and beyond the North Atlantic (Zhang et al., 2019; Ruprich-Robert et al., 2021, and references therein). However, the dominant driver of AMV remains highly contested; leading contenders include ocean dynamics (Arzel et al., 2022; Kim et al., 2018; Zhang et al., 2019), atmospheric dynamics (Cane et al., 2017; Clement et al., 2015), and variations in external forcing (L. N. Murphy et al., 2021; Klavans et al., 2022). Each of these drivers imply different timescales of predictability, and the short observational record further complicates the disentanglement of their contributions.

Yet the subpolar North Atlantic (SPNA), the center of action for AMV, is considered among the most predictable locations for SST and ocean heat content across all ocean basins, with skill extending to decadal timescales (Buckley et al., 2019; S. Yeager, 2020). Mean wintertime mixed-layer depths reach over 1,000 m within the SPNA, resulting in large heat capacity that translates to long persistence and memory of SST anomalies (Deser et al., 2003; Holte et al., 2017). The SPNA encompasses key deep-water formation sites of the Atlantic Meridional Overturning Circulation (AMOC), and has been linked to multi-year to multi-decadal predictability, both locally and in other regions such as the tropical Atlantic (Dunstone et al., 2011; Menary et al., 2015).

© 2023. The Authors.

This is an open access article under the terms of the [Creative Commons Attribution License](https://creativecommons.org/licenses/by/4.0/), which permits use, distribution and reproduction in any medium, provided the original work is properly cited.

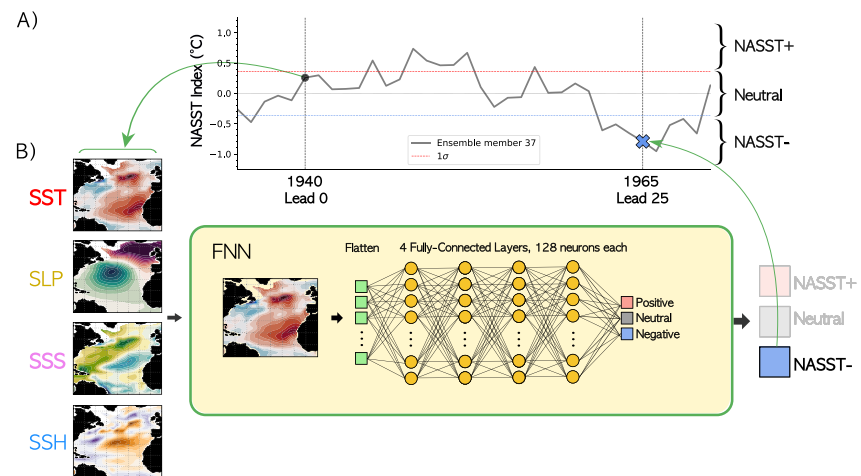


Figure 1. Schematic diagram of the NN prediction of NASST state using an example NASST– event in 1965 from ensemble member 37 of CESM1 LENS (Panel a). The snapshot of a selected predictor from 25 years prior (1940) is given to a FNN (Panel b), which outputs a prediction of the NASST state.

Current state-of-the-art approaches for decadal prediction of the climate system are often computationally intensive and highly sensitive to initial conditions, or constrained by assumptions of linearity in simplified models such as the Linear Inverse Model (Huddart et al., 2017; Meehl et al., 2022; Smith et al., 2019; Zanna, 2012). An alternative pathway emerges from neural networks (NN) and their ability to capture nonlinear processes and transformations (Hornik et al., 1989; Toms et al., 2020). NNs have successfully outperformed dynamical forecasts of El Niño–Southern Oscillation (ENSO) at interannual timescales (Ham et al., 2019) and detecting transitions between positive and negative states of the Pacific Decadal Oscillation (Gordon et al., 2021). Furthermore, recent developments of techniques such as Layer-wise Relevance Propagation (LRP) provide a way to peer into the “black box” of the NNs and identify the critical features for skillful predictions (Gordon et al., 2021; Toms et al., 2020; Wang et al., 2022). In this work, we investigate the potential of applying NNs to predicting North Atlantic sea surface temperatures (NASST) and use LRP to examine the relative importance of atmospheric and oceanic sources of predictability across multiple timescales.

2. Methods and Data

2.1. Data Sets

We use the Community Earth System Model 1 (CESM1) Large Ensemble Simulations (LENS) based on a fully-coupled global climate model with nominal 1-degree resolution (Kay et al., 2015). We focus on a single model to investigate if NNs can learn the physics of NASST variability, without confounding factors and biases that arise from cross-model comparisons. CESM1 LENS features 42 members under the same historical-era forcing from the Coupled-Model Intercomparison Project 5 (CMIP5), but with slightly different atmospheric initial conditions, representing a comprehensive range of intrinsic climate variability. We use the period common across all ensemble members (1920–2005), totaling of 3,612 years of data for training, validation, and testing of the NNs.

To investigate if the predictability learned from CESM1 translates to a realistic data set, we test the NNs on an observational data set, the Hadley Center Sea Ice and Sea Surface Temperature (HadISST) version 1 that includes monthly data between 1870 and 2022 at 1-degree resolution (Rayner et al., 2003). Since the NNs require inputs of the same size, we re-grid HadISST to match the CESM1 resolution using bilinear interpolation.

2.2. Prediction Objective

The input features are 2-D annual mean snapshots of atmospheric and/or oceanic predictors (discussed in Section 2.3) over the North Atlantic (80 to 0°W, 0 to 65°N), and the output prediction is the state of NASST (either positive, negative, or neutral) a given number of years later (Figure 1). The NASST index is the area-weighted,

annual mean SST anomaly over the North Atlantic, essentially the unfiltered AMV Index (Ting et al., 2009). Considering recent work that suggests the importance of external forcing in driving AMV (L. N. Murphy et al., 2021; Klavans et al., 2022), we also examine differences in predictability of NASST *with* and *without* external forcings such as the anthropogenic warming trend, defined by the 42-member ensemble mean (referred to as *forced* and *unforced*, respectively). This is performed prior to subsetting the data for training, validation, and testing.

We focus on predicting extreme NASST states due to its strong scientific and societal impacts. A 1-standard deviation (σ) threshold is used to separate the NASST into positive, negative, and neutral states (similar results are obtained using tercile thresholds). The threshold was selected to be high enough to distinguish extreme NASST anomalies, but low enough to permit sufficient samples for training. Framing this as a classification rather than a regression problem facilitates the application and interpretation of LRP output (Toms et al., 2020). To prevent biases toward predicting a specific class simply due to its frequency of occurrence, following standard practice (Buda et al., 2018; Drummond & Holte, 2003; Gordon et al., 2021), we subsample across CESM1 members set aside for training and validation so that there are equally 300 events per NASST state (see Section 2.4).

2.3. Atmospheric and Oceanic Predictors

To evaluate the importance of atmospheric versus oceanic drivers for NASST variability, we train networks to predict the NASST state given 2-D annual mean anomalies of the 4 following predictors:

1. *SST*, also used to calculate the NASST indices.
2. *Sea level pressure (SLP)*, an atmospheric predictor reflecting the state of the dominant atmospheric modes of variability in the region, for example, the North Atlantic Oscillation (NAO) (Hurrell & Deser, 2010; Ruprich-Robert & Cassou, 2015).
3. *Sea surface salinity (SSS)*, an oceanic predictor that is not directly damped by heat fluxes to the atmosphere, allowing for the investigation of redistribution and damping by ocean circulation and its connections with NASST variability (Zhang, 2017).
4. *Sea surface height (SSH)*, an oceanic predictor used to infer geostrophic circulation with connections to variations in the strength of subpolar gyre (Koul et al., 2020). SSH is also related to subsurface ocean heat content with potential for long-term predictability (Buckley et al., 2019; S. Yeager, 2020).

These predictors are observable from the ocean surface, and are thus more likely to have longer records into the future with satellite observations, providing potential for application to operational predictions of climate. We tested additional predictors from CESM1, including net air-sea heat flux, barotropic streamfunction, mixed-layer depth, heat and salt content, and wind stress and its curl. None of these predictors yielded significantly better performance, so we focus on the above four variables.

Each predictor is cropped to the domain used to compute the NASST index. Ocean variables are re-gridded to match atmospheric grid using bilinear interpolation. We exclude regions over land and where the ice fraction exceeds 5%. This allows us to compare oceanic and atmospheric predictors over shared areas where the signal is not dominated by sea ice variability, though including those points did not significantly impact the predictive skill. Each predictor is normalized to have a standard deviation of 1 across all dimensions, ensuring comparable variability between predictors and equal numerical contribution during the training process (Singh & Singh, 2020). Multiple NNs are trained with each of the above mentioned predictors separately. NNs that include all predictors as input did not yield improved skill, but rather indicate equivalent accuracy to the best predictor at each leadtime (not shown).

2.4. Network Architecture and Training Procedure

To separately investigate the dependency in timescale and predictor, each NN is trained to predict the NASST state at a specific leadtime ($t = 0$ to 25 years) given one predictor at a time. We withhold 10 members of CESM1 LENS for testing, and split the remaining 32 members into training (90%) and validation (10%) subsets. We initialize 100 different networks to account for randomness in the training process, totaling 10,400 networks (26 leadtimes \times 4 predictors \times 100 initialized networks). The training and validation sets are shuffled and resampled for each initialized network, ensuring that the results are not sensitive to a particular subset. Each network

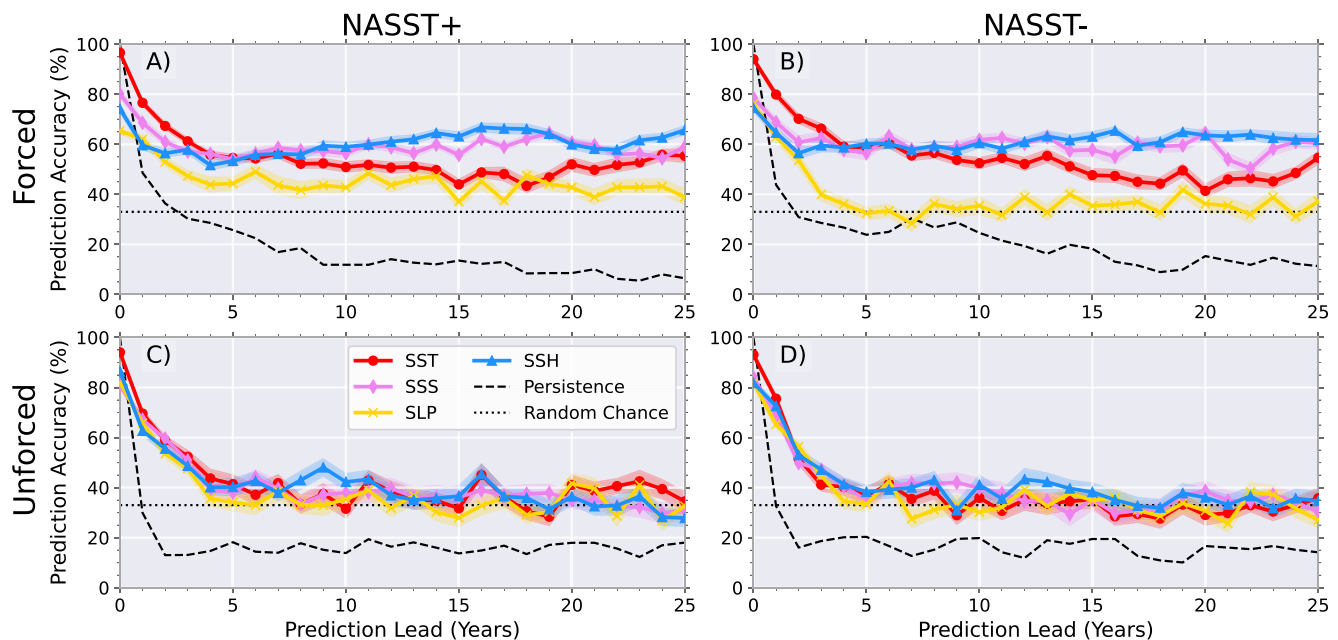


Figure 2. The mean accuracy by leadtime for predicting NASST+ and NASST− states for NNs trained with each predictor. X-axis is the prediction leadtime from 0 to 25 years. Shading indicates the 95% standard error of 100 NNs for each predictor. NNs trained with oceanic predictors SSH (blue) and SSS (pink) outperform those trained with SST (red) and SLP (yellow) at long leadtimes in the forced case (a–b). For the unforced case (c–d), performance is similar to the random chance baselines after 5–10 years (c–d).

is trained for 50 epochs, but the training process is stopped if the validation loss increases for five consecutive epochs to prevent over-fitting. All discussed results are from the same withheld testing set across all 100 networks.

We explored combinations of architectures and hyperparameters for convolutional neural networks (CNNs) and fully-connected neural networks (FNNs). FNNs have multiple layers consisting of individual neurons, each possessing trainable weights for all neurons in the preceding layer. In contrast, CNNs feature alternating convolutional layers, with trainable weights on filters that are applied throughout an image, and pooling layers to further reduce the dimensions, resulting in feature maps containing essential patterns needed for a prediction objective. Both architectures yielded comparable performance (Figure S4c in Supporting Information S1). Our hyperparameter tests revealed potential further optimization for SSH, but no systematic improvements across all predictors (see SI: Hyperparameter Testing). Since our objective is not to maximize accuracy, but rather to gain physical insight on drivers of NASST variability by examining inter-predictor differences, we focus on results for a simple 4-layer FNN with 128 neurons per layer and adopt the same architecture for all predictors.

2.5. Prediction Baselines

We compare the accuracy of the trained NNs to two baselines. Since each class is evenly sampled during the training, there is a 33% chance that a given class will occur, which we set as the *random chance baseline*. We additionally examine the other extreme using the standard *persistence baseline* that assumes uninterrupted continuation of the current state (A. H. Murphy, 1992). For example, if the system is at NASST+ at the starting time ($t = 0$ years), we assume it will also be NASST+ for the target leadtime.

3. Higher Skill From Oceanic Predictors at Multidecadal Leadtimes in the Presence of External Forcing

We focus on the prediction skill for NASST+ and NASST− events (Figure 2). For the predictions of Neutral events, the NNs had low accuracy equivalent to random chance. This is expected due to the challenge of predicting cases at the class boundaries or events with a weaker signal (Batista et al., 2004).

In the forced case (Figures 2a and 2b), NNs outperform both persistence and random chance baselines regardless of the predictor. The atmospheric variable, SLP, has similar-to-worse accuracy at all leadtimes compared to SST.

While this is unsurprising, considering the short persistence timescales of the atmosphere in the extratropics, on the order of weeks (Frankignoul & Hasselmann, 1977), the NN still outperforms the persistence and random chance baselines for predicting NASST+ across all leadtimes.

While SST is a better predictor at earlier leadtimes, NNs trained by both oceanic predictors (SSS and SSH) achieve consistently higher accuracy than SST at decadal and longer leadtimes (Figures 2a and 2b). Prolonged predictability from SSS could arise from absence of strong, direct damping by turbulent heat fluxes that exists in SSTs, allowing for more persistent SSS anomalies (Mignot & Frankignoul, 2003; Zhang, 2017). Similarly, subsurface heat content information present in SSH is shielded from damping by surface heat fluxes, leading to more persistence and potential predictability relative to SST (Buckley et al., 2019; Deser et al., 2003).

The increased predictability from oceanic variables is dependent upon the presence of external forcings. After removing the ensemble mean from the predictors and NASST index and repeating the training procedure, all NNs exhibit performance comparable to random chance after 5–10 years with minimal inter-predictor difference. Reduced damping of oceanic variables could lead to greater memory of externally forced signals in oceanic predictors. Overall, this highlights the importance of considering external forcing for climate prediction on multidecadal timescales and its enhancement of predictability derived from oceanic variables.

4. Consistent Source of Long-Term Predictability in the Transition Zone

To investigate the sources of predictability, we use LRP to examine the network's decision-making process (Bach et al., 2015; Böhle et al., 2019a). LRP back-propagates the “relevance” for given sample's prediction from the final output node to the input layer of the NN. The total relevance is conserved during this process through propagation rules, creating a “heatmap” of each pixel's contribution to the network's final decision (Montavon et al., 2019; Samek et al., 2021). We found that negative relevance values were highly sample and network dependent, and elected to use the $LRP_{\alpha\beta}$ rule (with $\alpha = 1.1$, $\beta = 0.1$, $p = 2$, $\epsilon = 10^{-2}$, see Section S3 in Supporting Information S1 for testing details), emphasizing positive contributions to a given sample that were consistent across samples (Binder et al., 2016). Previous works compared such relevance maps with known patterns of physical processes for predicting Pacific climate variability for possible correspondences (Gordon et al., 2021; Toms et al., 2020).

Since LRP produces the relevance map for a single sample, we examine the overall learned source of predictability by compositing relevances across *correct* predictions for the top 50 performing NNs of NASST+ and NASST−. The results are unchanged if all networks are included in the composite. The composites are normalized prior to visualization to have values between 0 and 1, though the raw output relevance is of order 10^{-4} . We show relevance composites for key leadtimes between 0 and 25 years overlaid on composites of input predictors at corresponding leadtimes (Figure 3) for the forced NASST+ cases. Results are broadly consistent in unforced and for NASST− cases (Figures S5 and S6 in Supporting Information S1).

For instantaneous predictions (leadtime 0), the relevance maps resemble known patterns associated with AMV and its drivers. For example, the SST relevance map (Figure 3e) captures the canonical horseshoe pattern of AMV (Zhang et al., 2019). Furthermore, the maximum relevance south of Newfoundland in SST, SSH, and SSS is collocated with the SPNA-Gulf Stream dipole associated with AMV-related SSTs and major ocean circulation features (Gu & Gervais, 2022; Nigam et al., 2018; Oelsmann et al., 2020; Zhang, 2008). Interestingly, a second relevance maxima for SSS is present near the Amazon River outflow region, though further investigation is needed to determine if this is a model-dependent feature and its physical mechanisms. Overall, these aspects lend confidence that the NN has learned to rely upon regions that vary strongly with AMV and its associated ocean drivers.

Patterns associated with atmospheric drivers of NASST variability also emerge in relevance maps at leadtimes longer than 5 years (Figures 3f–3i). Successful predictions by SLP-trained NNs rely upon negative SLP anomalies near the Icelandic Low in the northeastern Atlantic, a center of action for NAO (Deser et al., 2010; Hurrell & Deser, 2010). This learned reliance on the NAO-NASST linkage without additional input is encouraging, suggesting that additional predictability beyond the persistence baseline achieved by SLP-trained NNs may arise from large-scale air-sea interaction in this region and resulting ocean circulation anomalies.

The Transition Zone region between the subpolar and subtropical gyres is consistently important for predicting NASST regardless of leadtime for oceanic predictors (Figures 3k–3t) (Buckley & Marshall, 2016). This

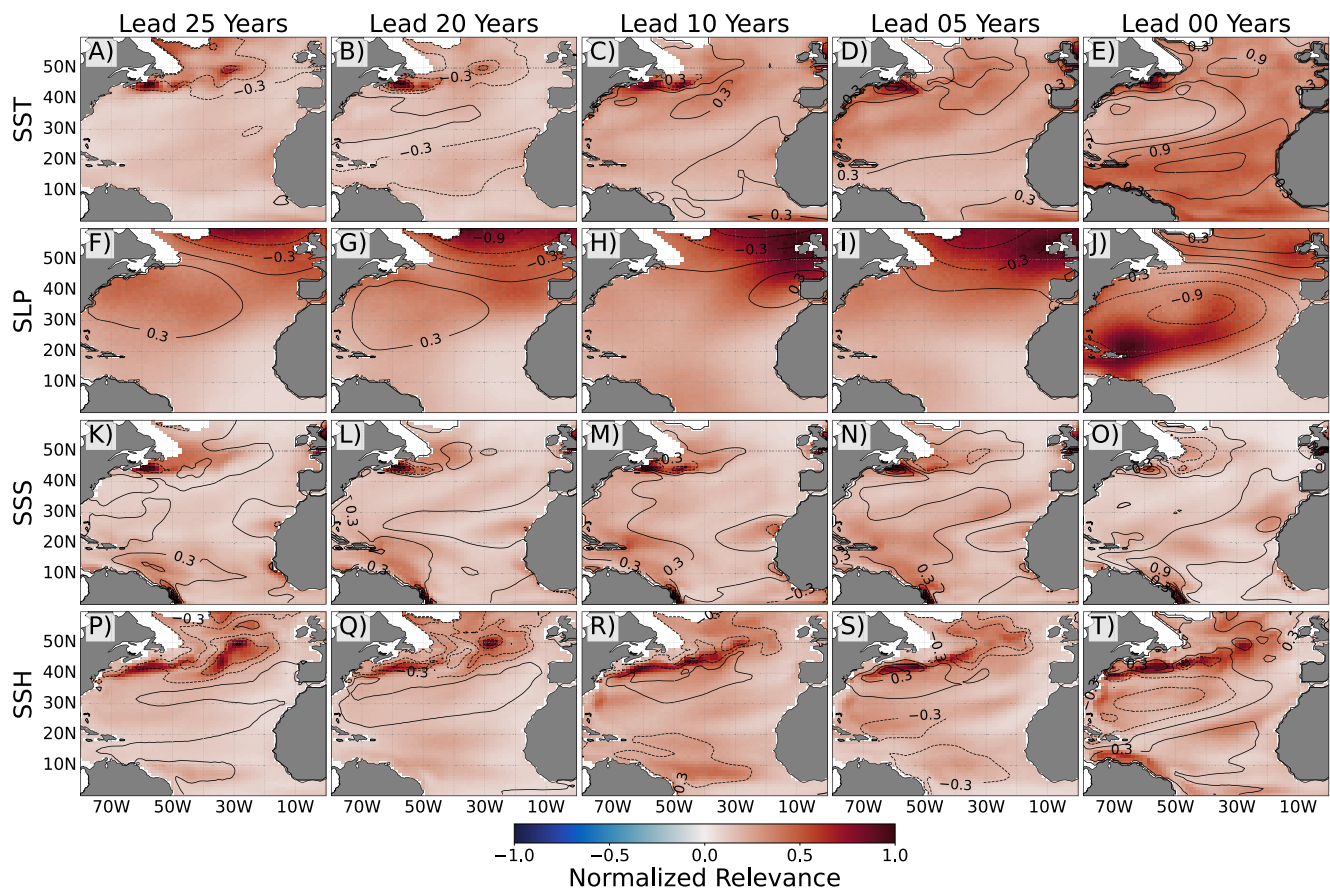


Figure 3. Composite relevance values (color) for “correct” NASST+ predictions of the top 50 performing networks for 0- to 25-year leadtimes, for the predictors from SST (a–e), SLP (f–j), SSS (k–o) and SSH (r–t), respectively. Relevance values are normalized for each composite. SSS relevance values were doubled to aid interpretability. Contours are the respective composites of standardized predictors for the given leadtime.

region has been connected to long-timescale oceanic processes, such as AMOC and its associated fingerprint in surface and subsurface temperatures (Zhang, 2008). In addition to prediction timescale, relevance over this region remains high irrespective of the class (NASST+ or NASST–) or the presence of external forcing (Figure S6 in Supporting Information S1). The potential importance of ocean dynamics for both forced and unforced NASST predictability is highlighted by this tendency for NNs to focus on this region.

5. CESM1-Trained Neural Networks Predict the Multidecadal Oscillation of Observed NASST States

Does the NNs' skill for NASST prediction apply beyond the CESM1 model world? Considering limited observational records of SSH, SSS, and SLP, we test if NNs trained on CESM1 SSTs can successfully predict the NASST state in HadISST. We pre-process and normalize the data using the same approach as for the CESM1 output (Section 2.3). Accounting for reductions due to the 25 years leadtime, there remains 128 years of data between 1895 and 2022. The 1σ threshold (0.55°C) yielded 29 (17) NASST+ (NASST–) events. The distribution is skewed due to the warming trend. Due to the limited samples, the accuracy values were noisy, particularly at long leadtimes. Therefore, we focus broadly on the frequency of predictions by class (Figure 4).

The frequency of predictions by class across all NNs aligns with multidecadal NASST oscillations in HadISST, with more frequent NASST– predictions coinciding with negative NASST index values pre-1925 and 1960–1990. This is true particularly for interannual and multidecadal leadtimes (Figures 4a and 4c), with shifted phasing at decadal leadtimes (Figure 4b). The same results are recovered for the unforced case, though the multidecadal phasing of predictions is nearly absent for the decadal leadtimes (Figure S7 in Supporting Information S1). These are surprising results for two main reasons: The first is that the NN is not simply predicting the anthropogenic

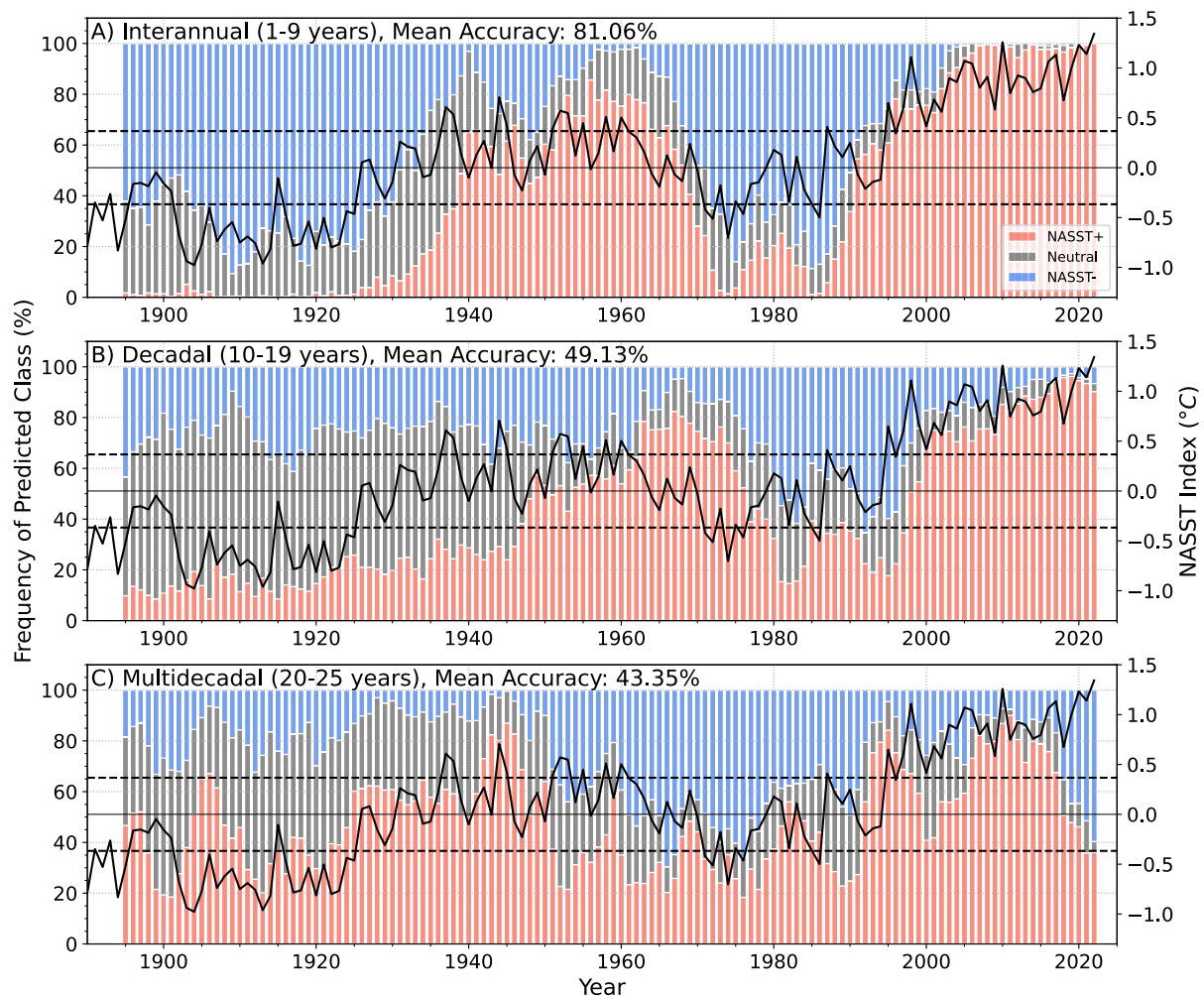


Figure 4. Frequency of predicted class of each target year aggregated for interannual (1–9 years) (a), decadal (10–19 years) (b), and multidecadal (20–25 years) (c) lead times for the HadISST (in colored bars) and corresponding mean accuracy across NASST+ and NASST– predictions (indicated in panel titles). Blue/red/gray bars are the frequency of the negative/positive/neutral NASST predictions. The NASST Index from HadISST (solid-black line) and 1σ thresholds (dashed-black lines) are shown for reference.

warming trend (e.g., monotonically increasing NASST+ predictions in time), but has instead successfully learned the non-linear, oscillatory behavior of the observed NASST index. The second is that the weights *not* have been re-adjusted to HadISST, revealing that NNs trained on potentially biased CESM1 output maintain their ability to predict the phasing of observed multidecadal climate variability. Overall, this joins a growing body of studies that suggests promise for applying NNs trained on model output to predicting the trajectory of non-linear multidecadal climate variability in corresponding observational data sets such as HadISST (Labe & Barnes, 2022).

6. Discussion and Summary

We investigated the potential of applying NNs to multidecadal prediction of NASST variability and used LRP to understand the contributions of oceanic and atmospheric drivers. Three main conclusions of this work are:

1. NNs trained with oceanic variables can predict NASST+ and NASST– states on multidecadal timescales, outperforming persistence and random chance baselines in the presence of external forcing.
2. The Transition Zone emerges as consistent region from where NNs derive predictive skill, regardless of prediction leadtime, NASST state, and the presence of the external forcing, suggesting a connection to ocean dynamics.
3. NNs trained on CESM1 were able to predict the multidecadal phasing of observed NASST states without weight readjustment, suggesting promise for training NNs using model output for multidecadal prediction of observed climate.

Previous studies have noted the importance of ocean initial conditions over external forcing for SST predictability over the SPNA (S. Yeager et al., 2018). Our results suggests that external forcing can enhance predictability derived from oceanic variables, evidenced by the difference in skill between atmospheric and oceanic predictors in the forced case at multidecadal timescales. A possible explanation is the larger heat capacity of the ocean allows for longer memory of externally forced signals, leading to enhanced predictability on multidecadal timescales (Frankignoul & Hasselmann, 1977). Further studies comparing regional variations in SST predictability and initial-state dependence, particularly with the SPNA, could yield further insight on the mechanisms behind increased skill of oceanic predictors under different forcing scenarios (S. Yeager et al., 2018; Gordon & Barnes, 2022).

A remarkably consistent feature across timescales in both unforced and forced cases is the high relevance over the Transition Zone. Decadal predictability from this region in CESM1 has been attributed to slow southward propagation of water mass anomalies from the North Atlantic Deep Water formation region (S. G. Yeager et al., 2015). While this southward communication of anomalies has led to the suggestion that Transition Zone variability is driven by AMOC-related processes, others argue that buoyancy anomalies originating in the Transition Zone are advected cyclonically around the subpolar gyre to the western boundary where they influence AMOC variations (Buckley & Marshall, 2016; Zhang, 2008). Regardless of the direction of causality and dynamic linkages between AMOC and Transition Zone anomalies, a common thread is the involvement of ocean dynamics for long-term NASST predictability (Little et al., 2020).

Predictability arising from a stationary feature over a region, rather than smaller-scale features propagating across the domain, could explain the comparable performance between FNNs and CNNs; For predicting NASST, the absolute position of the feature is more important than its translation invariance, erasing advantages conferred by the CNN's filters that specialize in capturing such features throughout the input (Barnes et al., 2022). Expanding the input domains could potentially reveal additional regions of predictability, particularly considering the recent interest in inter-basin interactions (Gordon & Barnes, 2022; Hong et al., 2022).

Our current approach uses LRP_{ap} , a method known to mix negative and positive relevances (Bommer et al., 2023; Mamalakis et al., 2022). Since this is more likely with increasing network complexity, we used a simpler 4-layer FNN (Mamalakis et al., 2022). Additionally, we find the negative relevances are largely inconsistent across samples and are reduced by the compositing operations. Investigations using additional explainability methods to assess the robustness of high relevance over the Transition Zone is critical future step (Bommer et al., 2023). Our hyperparameter testing indicated improved validation accuracy at long leadtimes for SSH-trained NNs with increased number of layers, suggesting pathways for further optimization (Figures S1–S3 in Supporting Information S1).

A cautionary note is that higher accuracy from networks trained with oceanic predictors could be a model dependent feature. Our results used CESM1, a coarse-resolution model with biases in the separation of the Gulf Stream and position of the North Atlantic Current (Kirtman et al., 2012). Since our relevance maps reveal that NNs depend upon this region for skillful predictions of NASST state, verifying the model dependence of this aspect by training NNs with other model large ensembles, reanalyses, or observational data sets is an important future endeavor. Considering connections between biases in mean state and decadal variability over the SPNA as well as sensitivity to external forcing, exploring correspondences between the resultant relevance maps and biases in ocean circulation may unveil further hints on the importance of ocean dynamics for NASST predictability (Menary et al., 2015).

Data Availability Statement

The monthly output from the CESM1 Large Ensemble is publicly available from the National Center for Atmospheric Research's Climate Data Gateway on the Earth System Grid (Kay et al., 2015; <https://www.cesm.ucar.edu/community-projects/lens/data-sets/>). Further specific instructions on accessing the CESM1 variables TS, LANDFRAC, ICEFRAC, SSS, PSL, and SSH used for this study CESM1 is detailed at this link (<https://www.cesm.ucar.edu/community-projects/lens/data-sets/>). The HadISST data set can be downloaded directly from their website (Rayner et al., 2003; <https://www.metoffice.gov.uk/hadobs/hadisst/>). Software for this work is available on Zenodo (Liu et al., 2023, DOI: <https://doi.org/10.5281/zenodo.8342739>), and the corresponding linked GitHub repository (https://github.com/glennliu265/predict_nasst). The Pytorch-LRP Software can be found in the following repository (Böhle et al., 2019b; <https://github.com/moboehle/Pytorch-LRP>).

Acknowledgments

GL is supported by the Department of Defense through the National Defense Science and Engineering Graduate Fellowship Program. GL and Y-OK gratefully acknowledge the support by the U.S. Department of Energy Office of Science Biological and Environmental Research as part of the Regional and Global Model Analysis program area (DE-SC0019492). Y-OK is also supported by National Science Foundation Division of Atmospheric and Geospace Sciences Climate and Large-scale Dynamics Program (AGS-2055236). PW acknowledges Grant 2128617 from the Atmospheric Chemistry Division of the National Science Foundation and support of VoLo foundation. The authors are very grateful to two anonymous reviewers and the editor Kristopher Karnauskas for their insightful and thorough comments.

References

- Arzel, O., Huck, T., Hochet, A., & Mussa, A. (2022). Internal ocean dynamics contribution to north Atlantic interdecadal variability strengthened by ocean–atmosphere thermal coupling. *Journal of Climate*, 35(24), 4605–4624. <https://doi.org/10.1175/jcli-d-22-0191.1>
- Bach, S., Binder, A., Montavon, G., Klauschen, F., Müller, K.-R., & Samek, W. (2015). On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. *PLoS One*, 10(7), e0130140. <https://doi.org/10.1371/journal.pone.0130140>
- Barnes, E. A., Barnes, R. J., Martin, Z. K., & Rader, J. K. (2022). This looks like that there: Interpretable neural networks for image tasks when location matters. *Artificial Intelligence for the Earth Systems*, 1(3), e220001. <https://doi.org/10.1175/ai-es-d-22-0001.1>
- Batista, G. E., Prati, R. C., & Monard, M. C. (2004). A study of the behavior of several methods for balancing machine learning training data. *ACM SIGKDD explorations newsletter*, 6(1), 20–29. <https://doi.org/10.1145/1007730.1007735>
- Binder, A., Montavon, G., Lapuschkin, S., Müller, K.-R., & Samek, W. (2016). Layer-wise relevance propagation for neural networks with local renormalization layers. In *Artificial neural networks and machine learning—ICANN 2016: 25th international conference on artificial neural networks, Barcelona, Spain, September 6–9, 2016, proceedings, part II* (Vol. 25, pp. 63–71).
- Böhle, M., Eitel, F., Weygandt, M., & Ritter, K. (2019a). Layer-wise relevance propagation for explaining deep neural network decisions in MRI-based Alzheimer's disease classification. *Frontiers in Aging Neuroscience*, 11, 194. <https://doi.org/10.3389/fnagi.2019.00194>
- Böhle, M., Eitel, F., Weygandt, M., & Ritter, K. (2019b). Pytorch-LRP: Basic LRP implementation in pytorch [Software]. GitHub. Retrieved from <https://github.com/mobochle/Pytorch-LRP>
- Bommer, P., Kretschmer, M., Hedström, A., Bareeva, D., & Höhne, M. M.-C. (2023). Finding the right XAI method—a guide for the evaluation and ranking of explainable AI methods in climate science. arXiv preprint arXiv:2303.00652.
- Buckley, M. W., DelSole, T., Lozier, M. S., & Li, L. (2019). Predictability of north Atlantic sea surface temperature and upper-ocean heat content. *Journal of Climate*, 32(10), 3005–3023. <https://doi.org/10.1175/jcli-d-18-0509.1>
- Buckley, M. W., & Marshall, J. (2016). Observations, inferences, and mechanisms of the Atlantic meridional overturning circulation: A review. *Reviews of Geophysics*, 54(1), 5–63. <https://doi.org/10.1002/2015rg000493>
- Buda, M., Maki, A., & Mazurowski, M. A. (2018). A systematic study of the class imbalance problem in convolutional neural networks. *Neural Networks*, 106, 249–259. <https://doi.org/10.1016/j.neunet.2018.07.011>
- Cane, M. A., Clement, A. C., Murphy, L. N., & Bellomo, K. (2017). Low-pass filtering, heat flux, and Atlantic multidecadal variability. *Journal of Climate*, 30(18), 7529–7553. <https://doi.org/10.1175/jcli-d-16-0810.1>
- Clement, A., Bellomo, K., Murphy, L. N., Cane, M. A., Mauritsen, T., Rädel, G., & Stevens, B. (2015). The Atlantic multidecadal oscillation without a role for ocean circulation. *Science*, 350(6258), 320–324. <https://doi.org/10.1126/science.1254390>
- Deser, C., Alexander, M. A., & Timlin, M. S. (2003). Understanding the persistence of sea surface temperature anomalies in midlatitudes. *Journal of Climate*, 16(1), 57–72. [https://doi.org/10.1175/1520-0442\(2003\)016<0057:utpos>2.0.co;2](https://doi.org/10.1175/1520-0442(2003)016<0057:utpos>2.0.co;2)
- Deser, C., Alexander, M. A., Xie, S.-P., & Phillips, A. S. (2010). Sea surface temperature variability: Patterns and mechanisms. *Annual Review of Marine Science*, 2(1), 115–143. <https://doi.org/10.1146/annurev-marine-120408-151453>
- Drummond, C., & Holte, R. C. (2003). C4.5, class imbalance, and cost sensitivity: Why under-sampling beats over-sampling. In *Workshop on learning from imbalanced datasets II* (Vol. 11(1–8)).
- Dunstone, N., Smith, D., & Eade, R. (2011). Multi-year predictability of the tropical Atlantic atmosphere driven by the high latitude north Atlantic ocean. *Geophysical Research Letters*, 38(14), L14701. <https://doi.org/10.1029/2011gl047949>
- Frankignoul, C., & Hasselmann, K. (1977). Stochastic climate models, Part II application to sea-surface temperature anomalies and thermocline variability. *Tellus*, 29(4), 289–305. <https://doi.org/10.3402/tellusa.v29i4.11362>
- Gordon, E. M., & Barnes, E. A. (2022). Incorporating uncertainty into a regression neural network enables identification of decadal state-dependent predictability in CESM2. *Geophysical Research Letters*, 49(15), e2022GL098635. <https://doi.org/10.1029/2022gl098635>
- Gordon, E. M., Barnes, E. A., & Hurrell, J. W. (2021). Oceanic harbingers of pacific decadal oscillation predictability in CESM2 detected by neural networks. *Geophysical Research Letters*, 48(21), e2021GL095392. <https://doi.org/10.1029/2021gl095392>
- Gu, Q., & Gervais, M. (2022). Diagnosing two-way coupling in decadal north Atlantic SST variability using time-evolving self-organizing maps. *Geophysical Research Letters*, 49(8), e2021GL096560. <https://doi.org/10.1029/2021gl096560>
- Ham, Y.-G., Kim, J.-H., & Luo, J.-J. (2019). Deep learning for multi-year ENSO forecasts. *Nature*, 573(7775), 568–572. <https://doi.org/10.1038/s41586-019-1559-7>
- Holte, J., Talley, L. D., Gilson, J., & Roemmich, D. (2017). An Argo mixed layer climatology and database. *Geophysical Research Letters*, 44(11), 5618–5626. <https://doi.org/10.1002/2017gl073426>
- Hong, J.-S., Yeh, S.-W., & Yang, Y.-M. (2022). Interbasin interactions between the pacific and Atlantic oceans depending on the phase of pacific decadal oscillation and Atlantic multidecadal oscillation. *Journal of Climate*, 35(9), 2883–2894. <https://doi.org/10.1175/jcli-d-21-0408.1>
- Hornik, K., Stinchcombe, M., & White, H. (1989). Multilayer feedforward networks are universal approximators. *Neural Networks*, 2(5), 359–366. [https://doi.org/10.1016/0893-6080\(89\)90020-8](https://doi.org/10.1016/0893-6080(89)90020-8)
- Huddart, B., Subramanian, A., Zanna, L., & Palmer, T. (2017). Seasonal and decadal forecasts of Atlantic sea surface temperatures using a linear inverse model. *Climate Dynamics*, 49(5–6), 1833–1845. <https://doi.org/10.1007/s00382-016-3375-1>
- Hurrell, J. W., & Deser, C. (2010). North Atlantic climate variability: The role of the North Atlantic oscillation. *Journal of Marine Systems*, 79(3–4), 231–244. <https://doi.org/10.1016/j.jmarsys.2009.11.002>
- Kay, J. E., Deser, C., Phillips, A., Mai, A., Hannay, C., Strand, G., et al. (2015). The community earth system model (CESM) large ensemble project: A community resource for studying climate change in the presence of internal climate variability [Dataset]. Bulletin of the American Meteorological Society, 96(8), 1333–1349. <https://doi.org/10.1175/bams-d-13-00255.1>
- Kim, W. M., Yeager, S. G., & Danabasoglu, G. (2018). Key role of internal ocean dynamics in Atlantic multidecadal variability during the last half century. *Geophysical Research Letters*, 45(24), 13–449. <https://doi.org/10.1029/2018gl080474>
- Kirtman, B. P., Bitz, C., Bryan, F., Collins, W., Dennis, J., Hearn, N., et al. (2012). Impact of ocean model resolution on CCSM climate simulations. *Climate Dynamics*, 39(6), 1303–1328. <https://doi.org/10.1007/s00382-012-1500-3>
- Klavans, J. M., Clement, A. C., Cane, M. A., & Murphy, L. N. (2022). The evolving role of external forcing in north Atlantic SST variability over the last millennium. *Journal of Climate*, 35(9), 2741–2754. <https://doi.org/10.1175/jcli-d-21-0338.1>
- Koul, V., Tesdal, J.-E., Bersch, M., Hátún, H., Brune, S., Borchert, L., et al. (2020). Unraveling the choice of the north Atlantic subpolar gyre index. *Scientific Reports*, 10(1), 1–12. <https://doi.org/10.1038/s41598-020-57790-5>
- Labe, Z. M., & Barnes, E. A. (2022). Predicting slowdowns in decadal climate warming trends with explainable neural networks. *Geophysical Research Letters*, 49(9), e2022GL098173. <https://doi.org/10.1029/2022gl098173>
- Little, C. M., Zhao, M., & Buckley, M. W. (2020). Do surface temperature indices reflect centennial-timescale trends in Atlantic meridional overturning circulation strength? *Geophysical Research Letters*, 47(22), e2020GL090888. <https://doi.org/10.1029/2020gl090888>

- Liu, G., Wang, P., & Kwon, Y.-O. (2023). glennliu265/predict_nasst: Revision_update (v.1.1.0) [Software]. Zenodo. <https://doi.org/10.5281/zenodo.10161304>
- Mamalakis, A., Ebert-Uphoff, I., & Barnes, E. A. (2022). Neural network attribution methods for problems in geoscience: A novel synthetic benchmark dataset. *Environmental Data Science*, 1, e8. <https://doi.org/10.1017/eds.2022.7>
- Meehl, G. A., Teng, H., Smith, D., Yeager, S., Merryfield, W., Doblas-Reyes, F., & Glanville, A. A. (2022). The effects of bias, drift, and trends in calculating anomalies for evaluating skill of seasonal-to-decadal initialized climate predictions. *Climate Dynamics*, 59(11–12), 3373–3389. <https://doi.org/10.1007/s00382-022-06272-7>
- Menary, M. B., Hodson, D. L., Robson, J. I., Sutton, R. T., Wood, R. A., & Hunt, J. A. (2015). Exploring the impact of cmip5 model biases on the simulation of north Atlantic decadal variability. *Geophysical Research Letters*, 42(14), 5926–5934. <https://doi.org/10.1002/2015gl064360>
- Mignot, J., & Frankignoul, C. (2003). On the interannual variability of surface salinity in the atlantic. *Climate Dynamics*, 20(6), 555–565. <https://doi.org/10.1007/s00382-002-0294-0>
- Montavon, G., Binder, A., Lapuschkin, S., Samek, W., & Müller, K.-R. (2019). Layer-wise relevance propagation: An overview. *Explainable AI: Interpreting, explaining and visualizing deep learning*, 193–209. https://doi.org/10.1007/978-3-030-28954-6_10
- Murphy, A. H. (1992). Climatology, persistence, and their linear combination as standards of reference in skill scores. *Weather and Forecasting*, 7(4), 692–698. [https://doi.org/10.1175/1520-0434\(1992\)007<0692:cpatlc>2.0.co;2](https://doi.org/10.1175/1520-0434(1992)007<0692:cpatlc>2.0.co;2)
- Murphy, L. N., Klavans, J. M., Clement, A. C., & Cane, M. A. (2021). Investigating the roles of external forcing and ocean circulation on the Atlantic multidecadal SST variability in a large ensemble climate model hierarchy. *Journal of Climate*, 34(12), 4835–4849. <https://doi.org/10.1175/jcli-d-20-0167.1>
- Nigam, S., Ruiz-Barradas, A., & Chafik, L. (2018). Gulf stream excursions and sectional detachments generate the decadal pulses in the Atlantic multidecadal oscillation. *Journal of Climate*, 31(7), 2853–2870. <https://doi.org/10.1175/jcli-d-17-0010.1>
- Oelsmann, J., Borchert, L., Hand, R., Baehr, J., & Jungclaus, J. H. (2020). Linking ocean forcing and atmospheric interactions to Atlantic multidecadal variability in MPI-ESM1. 2. *Geophysical Research Letters*, 47(10), e2020GL087259. <https://doi.org/10.1029/2020gl087259>
- Rayner, N., Parker, D. E., Horton, E., Folland, C. K., Alexander, L. V., Rowell, D., et al. (2003). Global analyses of sea surface temperature, sea ice, and night marine air temperature since the late nineteenth century [Dataset]. *Journal of Geophysical Research*, 108(D14), 4407. <https://doi.org/10.1029/2002jd002670>
- Ruprich-Robert, Y., & Cassou, C. (2015). Combined influences of seasonal east Atlantic pattern and north Atlantic oscillation to excite Atlantic multidecadal variability in a climate model. *Climate Dynamics*, 44(1–2), 229–253. <https://doi.org/10.1007/s00382-014-2176-7>
- Ruprich-Robert, Y., Moreno-Chamarro, E., Levine, X., Bellucci, A., Cassou, C., Castruccio, F., et al. (2021). Impacts of Atlantic multidecadal variability on the tropical pacific: A multi-model study. *Npj Climate and Atmospheric Science*, 4(1), 33. <https://doi.org/10.1038/s41612-021-00188-5>
- Samek, W., Montavon, G., Lapuschkin, S., Anders, C. J., & Müller, K.-R. (2021). Explaining deep neural networks and beyond: A review of methods and applications. *Proceedings of the IEEE*, 109(3), 247–278. <https://doi.org/10.1109/jproc.2021.3060483>
- Singh, D., & Singh, B. (2020). Investigating the impact of data normalization on classification performance. *Applied Soft Computing*, 97, 105524. <https://doi.org/10.1016/j.asoc.2019.105524>
- Smith, D., Eade, R., Scaife, A., Caron, L.-P., Danabasoglu, G., DelSole, T., et al. (2019). Robust skill of decadal climate predictions. *Npj Climate and Atmospheric Science*, 2(1), 13. <https://doi.org/10.1038/s41612-019-0071-y>
- Ting, M., Kushnir, Y., Seager, R., & Li, C. (2009). Forced and internal twentieth-century SST trends in the north Atlantic. *Journal of Climate*, 22(6), 1469–1481. <https://doi.org/10.1175/2008jcli2561.1>
- Toms, B. A., Barnes, E. A., & Ebert-Uphoff, I. (2020). Physically interpretable neural networks for the geosciences: Applications to earth system variability. *Journal of Advances in Modeling Earth Systems*, 12(9), e2019MS002002. <https://doi.org/10.1029/2019ms002002>
- Wang, P., Yuval, J., & O’Gorman, P. A. (2022). Non-local parameterization of atmospheric subgrid processes with neural networks. *Journal of Advances in Modeling Earth Systems*, 14(10), e2022MS002984. <https://doi.org/10.1029/2022ms002984>
- Yeager, S. (2020). The abyssal origins of north Atlantic decadal predictability. *Climate Dynamics*, 55(7–8), 2253–2271. <https://doi.org/10.1007/s00382-020-05382-4>
- Yeager, S., Danabasoglu, G., Rosenbloom, N., Strand, W., Bates, S., Meehl, G., et al. (2018). Predicting near-term changes in the earth system: A large ensemble of initialized decadal prediction simulations using the community earth system model. *Bulletin of the American Meteorological Society*, 99(9), 1867–1886. <https://doi.org/10.1175/bams-d-17-0098.1>
- Yeager, S. G., Karspeck, A. R., & Danabasoglu, G. (2015). Predicted slowdown in the rate of Atlantic sea ice loss. *Geophysical Research Letters*, 42(24), 10–704. <https://doi.org/10.1002/2015gl065364>
- Zanna, L. (2012). Forecast skill and predictability of observed Atlantic sea surface temperatures. *Journal of Climate*, 25(14), 5047–5056. <https://doi.org/10.1175/jcli-d-11-00539.1>
- Zhang, R. (2008). Coherent surface-subsurface fingerprint of the Atlantic meridional overturning circulation. *Geophysical Research Letters*, 35(20), L20705. <https://doi.org/10.1029/2008gl035463>
- Zhang, R. (2017). On the persistence and coherence of subpolar sea surface temperature and salinity anomalies associated with the Atlantic multidecadal variability. *Geophysical Research Letters*, 44(15), 7865–7875. <https://doi.org/10.1002/2017gl074342>
- Zhang, R., Sutton, R., Danabasoglu, G., Kwon, Y.-O., Marsh, R., Yeager, S. G., et al. (2019). A review of the role of the Atlantic meridional overturning circulation in Atlantic multidecadal variability and associated climate impacts. *Reviews of Geophysics*, 57(2), 316–375. <https://doi.org/10.1029/2019rg000644>