



## RESEARCH ARTICLE

10.1029/2024JH000243

# Using Explainable AI and Transfer Learning to Understand and Predict the Maintenance of Atlantic Blocking With Limited Observational Data

 Huan Zhang<sup>1</sup>, Justin Finkel<sup>2</sup> , Dorian S. Abbot<sup>3</sup> , Edwin P. Gerber<sup>1</sup>, and Jonathan Weare<sup>1</sup> 
<sup>1</sup>Courant Institute of Mathematical Sciences, New York University, New York, NY, USA, <sup>2</sup>Department of Earth, Atmospheric, and Planetary Sciences, Massachusetts Institute of Technology, Cambridge, MA, USA, <sup>3</sup>Department of the Geophysical Sciences, University of Chicago, Chicago, IL, USA

## Key Points:

- Given sufficient training data, convolutional neural networks can predict the maintenance of Atlantic blocking from an initial blocked state
- Transfer learning from an idealized model to reanalysis data enables predictive skill in the low data regime of the observational record
- Feature importance analysis reveals the influence of upstream flow on blocking persistence and quantifies biases in the idealized model

## Correspondence to:

 H. Zhang,  
hz1994@nyu.edu

## Citation:

 Zhang, H., Finkel, J., Abbot, D. S., Gerber, E. P., & Weare, J. (2024). Using explainable AI and transfer learning to understand and predict the maintenance of Atlantic blocking with limited observational data. *Journal of Geophysical Research: Machine Learning and Computation*, 1, e2024JH000243. <https://doi.org/10.1029/2024JH000243>

Received 12 APR 2024

Accepted 21 OCT 2024

## Author Contributions:

**Conceptualization:** Huan Zhang, Justin Finkel, Dorian S. Abbot, Edwin P. Gerber, Jonathan Weare  
**Investigation:** Huan Zhang, Justin Finkel, Dorian S. Abbot, Edwin P. Gerber, Jonathan Weare  
**Methodology:** Huan Zhang, Justin Finkel, Dorian S. Abbot, Edwin P. Gerber, Jonathan Weare  
**Software:** Huan Zhang, Justin Finkel  
**Validation:** Huan Zhang, Justin Finkel  
**Visualization:** Huan Zhang, Justin Finkel  
**Writing – original draft:** Huan Zhang, Justin Finkel

**Abstract** Blocking events are an important cause of extreme weather, especially long-lasting blocking events that trap weather systems in place. The duration of blocking events is, however, underestimated in climate models. Explainable Artificial Intelligence are a class of data analysis methods that can help identify physical causes of prolonged blocking events and diagnose model deficiencies. We demonstrate this approach on an idealized quasigeostrophic (QG) model developed by Marshall and Molteni (1993), [https://doi.org/10.1175/1520-0469\(1993\)050<1792:taduop>2.0.co;2](https://doi.org/10.1175/1520-0469(1993)050<1792:taduop>2.0.co;2). We train a convolutional neural network (CNN), and subsequently, build a sparse predictive model for the persistence of Atlantic blocking, conditioned on an initial high-pressure anomaly. Shapley Additive ExPlanation (SHAP) analysis reveals that high-pressure anomalies in the American Southeast and North Atlantic, separated by a trough over Atlantic Canada, contribute significantly to prediction of sustained blocking events in the Atlantic region. This agrees with previous work that identified precursors in the same regions via wave train analysis. When we apply the same CNN to blockings in the ERA5 atmospheric reanalysis, there is insufficient data to accurately predict persistent blocks. We partially overcome this limitation by pre-training the CNN on the plentiful data of the Marshall-Molteni model, and then using Transfer learning (TL) to achieve better predictions than direct training. SHAP analysis before and after TL allows a comparison between the predictive features in the reanalysis and the QG model, quantifying dynamical biases in the idealized model. This work demonstrates the potential for machine learning methods to extract meaningful precursors of extreme weather events and achieve better prediction using limited observational data.

**Plain Language Summary** Blocking events are an important cause of extreme weather, especially long-lasting blocking events that trap weather systems in place. The duration of blocking events is, however, systematically underestimated in climate models. Using data generated by a simplified atmospheric model we demonstrate that, given sufficient training data, convolutional neural networks can predict the maintenance of Atlantic blocking from an initial blocked state. Next, we show that first training the neural network on data from the simplified model and then fine tuning the training using real world weather data enables prediction even with few examples of long-lasting blocking events in the observational record. Subsequent feature analysis of the resulting neural networks identifies the input variables that most strongly impact their predictions, revealing that areas of high pressure in certain parts of North America and the North Atlantic Ocean are important for predicting long-lasting blocking events and quantifying biases in the idealized model relative to real weather.

## 1. Introduction

Blocking events are high-amplitude, quasi-stationary anticyclonic high-pressure anomalies that give rise to prolonged abnormal weather conditions in the mid-to-high latitudes (Lupo, 2021; Rex, 1950; Woollings et al., 2018). Blocking events can lead to regional extreme weather by disrupting the usual westerly flow for extended periods (e.g., Kautz et al., 2022), causing extreme heatwaves, floods, and winter storms (e.g., Lupo et al., 2012).

The predictive skill of numerical weather models has improved dramatically, but they still cannot accurately forecast important aspects of blocking events. Blocking frequency and duration are generally simulated poorly by climate models (Davini & D'Andrea, 2020), and even by numerical weather prediction models in medium-range forecasts (Ferranti et al., 2015; Matsueda, 2009; Woollings et al., 2018). Several possible contributing factors have been proposed, including the accuracy of the model's mean flow (Scaife et al., 2010) or synoptic eddies

© 2024 The Author(s). Journal of Geophysical Research: Machine Learning and Computation published by Wiley Periodicals LLC on behalf of American Geophysical Union.

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.

Writing – review & editing: Huan Zhang,  
Justin Finkel, Dorian S. Abbot, Edwin  
P. Gerber, Jonathan Weare

(Berckmans et al., 2013; Zappa et al., 2014a), the model's resolution (Davini & D'Andrea, 2016) and subgrid-scale parameterizations (d'Andrea et al., 1998), and even the choice of blocking index itself (Dole & Gordon, 1983; Pelly & Hoskins, 2003; Tibaldi & Molteni, 1990).

Two commonly used blocking indices (Dole & Gordon, 1983; Tibaldi & Molteni, 1990) highlight two essential features of a blocking *event*: (a) a large positive anomaly of geopotential height that displaces the midlatitude jet, “blocking” the flow, that (b) persists for longer than typical synoptic variability. Often a five-day (5d) threshold is invoked, but the longer the flow remains in a blocked state, the more severe the implications, either for extended cold/hot conditions or an increased likelihood of compound storm events (i.e., back-to-back storms, which can dramatically increase the potential for damage; Kautz et al., 2022). The persistence of blocking is the focus of our study: given the onset of a blocked state, what is the likelihood that the flow will remain blocked for an extended period, 5 days for a standard event, or up to 9 days for more extreme cases? We take a data-driven approach, training a CNN to identify persistent blocks at the onset of a blocked state.

To understand blocking, various low-order models have been formulated to identify essential features. In an influential early work, Charney and DeVore (1979) modeled blocking as one of two equilibrium states of a set of dynamical equations for a highly truncated barotropic channel model. Others used low-order models to propose that the positive feedback of synoptic-scale eddies on the blocking structure contributes to the long-time maintenance of blocks (Hoskins et al., 1983; McWilliams, 1980; Shutts, 1983). While these low-order models have provided useful physical insight, their application to the real world is limited by lack of land-sea interactions, topography, and other factors. Comprehensive models, on the other hand, are becoming skillful in simulating realistic blocking (e.g., Davini et al., 2021), but their complexity makes it challenging to isolate the essential mechanism(s), and expensive to simulate numerous events.

To strike a balance between complexity, transparency, and statistical robustness from abundant data (model output), we begin with the Marshall-Molteni (MM) model (Marshall & Molteni, 1993), a three-layer quasi-geostrophic (QG) approximation of the atmosphere that has previously been used to study blocking events (e.g., Lucarini & Gritsun, 2020). The MM model captures the main features of the northern hemisphere atmosphere reasonably well. For example, Michelangeli and Vautard (1998) found that an enhanced baroclinic wavetrain traveling across the North Atlantic is necessary to trigger the onset of the Euro-Atlantic blocking in both this simple model and reanalysis. They also pointed out that wave-wave interactions and wave-mean interactions dominate local amplification and the propagation of anomalies, respectively.

The MM model allows us the freedom to develop and test methods in a data-rich setting, and precisely quantify the degradation of skill as we pass to a more realistic, data-poor setting. For the particular application of blocking, here we address the question: how well can a data-driven method identify persistent events as a function of the input data you allow it? Furthermore, to gain insight into the physics and predictability of blocking, we turn to Explainable Artificial Intelligence (XAI) techniques, following work by Labe and Barnes (2021) and Rampal et al. (2022). Specifically, we employ Shapley Additive ExPlanation (SHAP) analysis to identify key regions upstream of the blocking center that enable prediction, and use this to construct low-order models that can be interpreted in the context of prior work.

Our ultimate goal, however, is to forecast and understand the maintenance of blocks in our atmosphere, for which we shift the focus to ERA5 reanalysis (Hersbach et al., 2020). For the most extreme case of a 9-day block in the North Atlantic, only 18 have occurred in the historical record (See Table 3). What chance does a data-driven approach have? To address the problem of limited data, we apply TL: first we train a CNN on the MM model to learn the basic features of blocking, and then we re-train it on the limited ERA5 data to calibrate it for the real atmosphere. In this direction our results serve as proof-of-concept. It is likely another choice of physical model could strike a better balance between accuracy and simulation cost for our purpose. Nonetheless, we find that pre-training on the MM model yields a better predictor than when we train the same network on ERA5 alone, proving the efficacy of the TL approach.

The remainder of this paper is organized as follows. Section 2 introduces the Marshall-Molteni (MM) model, training data and blocking index. Section 3 formulates the blocking event criteria and forecasting problem. Section 4 discusses our CNN structure and training details. We first focus exclusively on the MM model in Sections 5 and 6, applying XAI techniques to visualize the important features for prediction and testing the results by building a sparse model with features guided by the XAI. We also suggest physical interpretations for these

**Table 1**  
Length of Trajectory (in Thousands of Days) Versus Number of Nascent Blocking States ( $T = 1$ ) in Training Set and Test Sets of Varying Size

Training data		Test data	
Days	Nascent blocked states	Days	Nascent blocked states
1k	63	250k	17,755
10k	699		
100k	7,024		
500k	35,078		
1000k	70,635		

predictive features. Finally, we turn to the ERA5 data set in Section 7, applying TL to improve the prediction of persistent blocks in ERA5, especially for more extreme events. SHAP analysis shows how TL has modified the CNN to adapt to the new data set, but preserves the use of key upstream regions for prediction.

## 2. Model and Blocking Index

Marshall and Molteni (1993) developed a 3-layer quasi-geostrophic model of the atmosphere to study atmospheric low-frequency variability. We refer the reader to Appendix A for a complete description. We use a Northern Hemisphere only version of the model developed by Lucarini and Gritsun (2020) with 6,210 degrees of freedom. The model is run with T31 horizontal resolution (corresponding to  $90$  longitude  $\times$   $23$  latitude gridpoints across the northern hemisphere). All model output fields, as well as the reanalysis used later, are averaged daily.

We use an index developed by Dole and Gordon (1983) to define blocking events, hereafter referred to as the DG index. This is an anomaly based blocking index, but has been shown to capture the same essential features of blocking as other measures, for example, that of Tibaldi and Molteni (1990). We compute this index by transforming the spherical harmonic representation of the streamfunction  $\psi$  at 500 hPa into approximate geopotential height,  $Z$ , on a Gaussian grid for latitude and a uniform grid for longitude. The approximation is the choice of a fixed Coriolis parameter  $f_0$  to convert from  $\psi$  to  $Z$ , which causes minimal distortion over our midlatitude area of focus. Blocks are based on deviations of the geopotential height from climatology, denoted  $Z'$ .

A *blocking event* is said to occur at a specific location when  $Z'$  stays above a tunable geopotential height anomaly threshold,  $M$ , for at least five consecutive days. In their paper, Dole and Gordon (1983) tested statistics for varying  $M$  values, ranging from 50 to 250 m, with subsequent studies adopting different thresholds (Chan et al., 2019; Table 2). For our investigation, we calibrated  $M = 100$  m for our MM model simulation to roughly match the blocking fraction computed from ERA5 reanalysis data, where we used the threshold  $M = 150$  m as in Mullen (1987).

Figure 1 shows the blocking event statistics during the simulation. For comparison, blocking event statistics computed from ERA5 reanalysis data from 1959 to 2021 are also shown. In this study, we focus on North Atlantic blockings indicated by the white rectangle in Figure 1. We pick this region because it has a relatively high blocking frequency, and for its important influence on western Europe. We use  $Z_B$ , the average 500 hPa geopotential height anomaly over this target region over the North Atlantic, to define blocked states and blocking events.

## 3. Probabilistic Forecasting and Event Definition

We aim to study the *maintenance* of blocks rather than their *onset*. Precisely, we formulate the question as the classification problem posed in Figure 2: given a nascent blocked state, that is, the state on a day that geopotential height anomalies over the North Atlantic first exceed the threshold  $M$ , can we immediately predict whether the flow will remain blocked for 5 or more days—evolving into a *blocking event*—or will the flow return back toward the climatological state before 5 days have passed? In the MM model, nascent blocked states evolve into 5-day persistent blocking events approximately 21% of the time.

We pose the classification problem: given only the state at the time of blocking onset, can a data-driven method accurately identify the rarer cases that will persist for more than 5 consecutive days? Mathematically, we denote the full model state by  $X$  and further introduce a variable  $T$  for the running duration of a blocked state:

$$T = (\text{days since } Z_B < M). \quad (1)$$

**Table 2**  
The Statistics of Blocking Events in Our MM 1,250k Day Simulation

Threshold	$Y = 1$	$Y = 0$	Positive rate
$\geq 5$ d	18,748	69,642	0.212
$\geq 7$ d	8,522	79,868	0.096
$\geq 9$ d	3,891	84,499	0.044

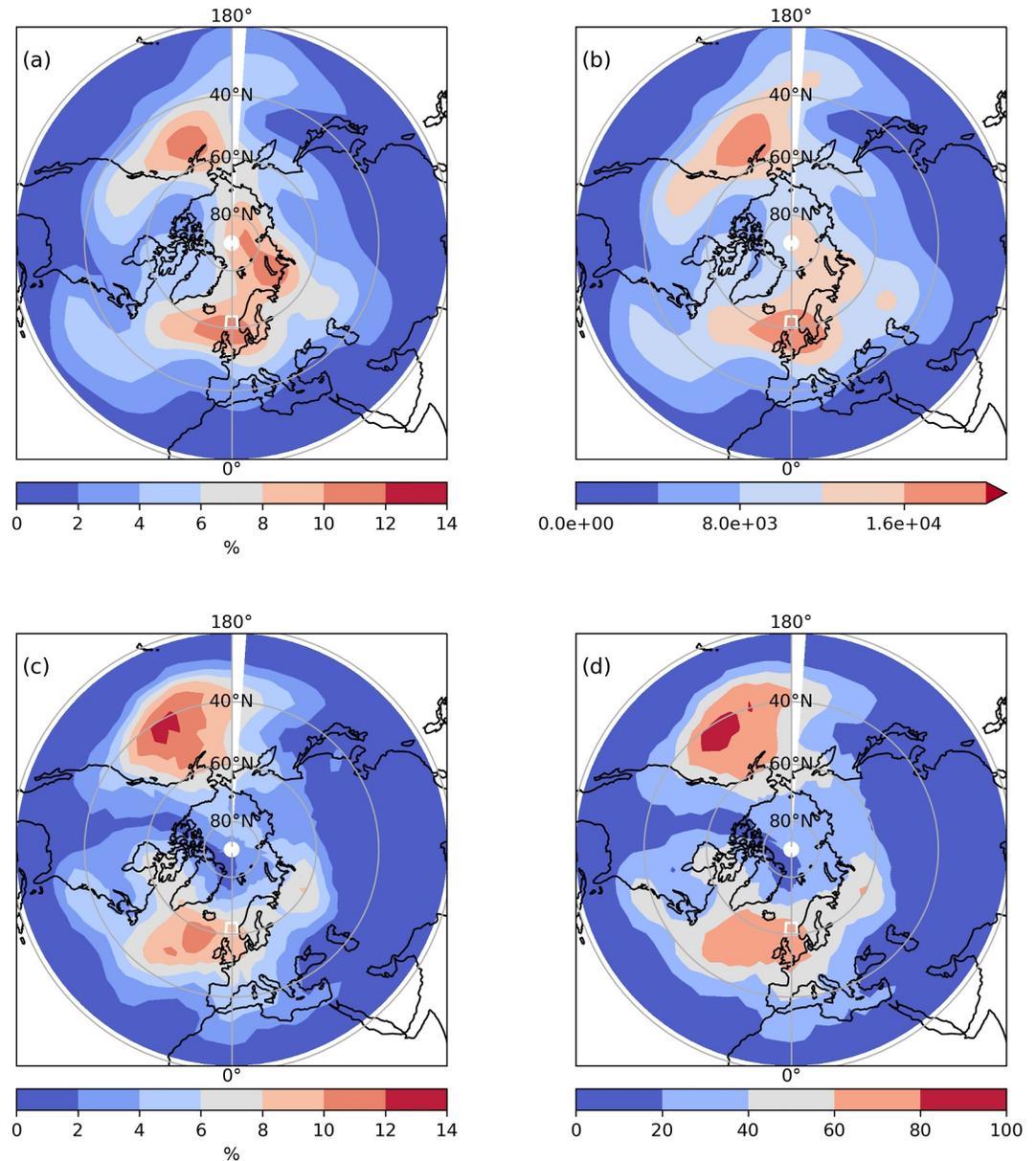
*Note.* The full data set exhibits 88,390 nascent blocking states ( $T = 1$  states).  $Y = 1$  marks the number of these nascent blocks that persist for 5, 7, or 9 days, thus evolving into a blocking event under these respective thresholds, while  $Y = 0$  denotes the number that don't make it to the threshold.

**Table 3**  
The Statistics of ERA5 Data Set in 1940–2022 December, January and February With  $T = 1$

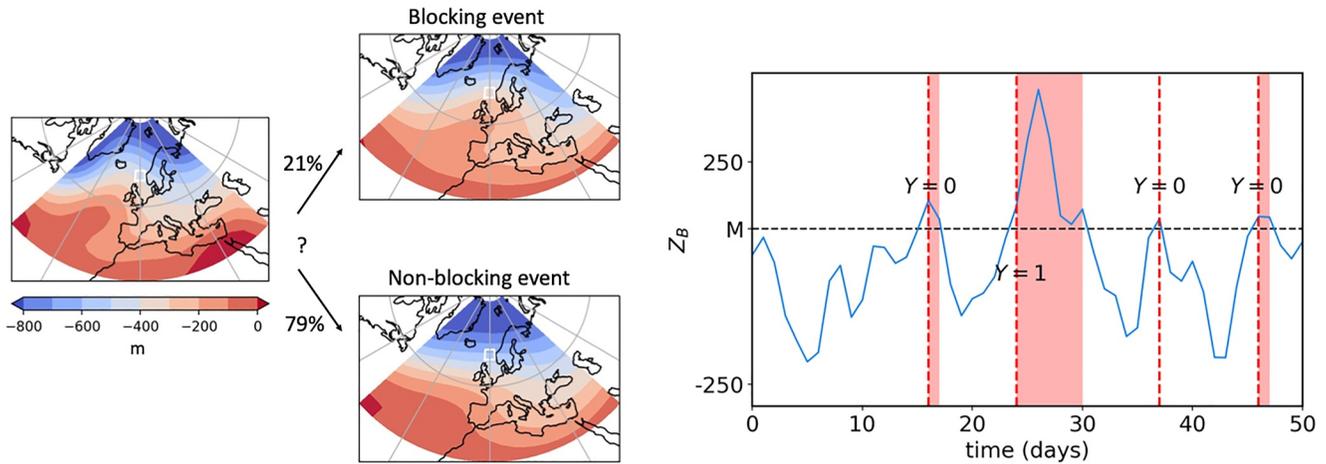
Threshold	$Y = 1$	$Y = 0$
$\geq 5$ d	84	189
$\geq 7$ d	36	237
$\geq 9$ d	18	255

Note that  $Z_B(t)$  is determined by the state vector  $X(t)$  at any time  $t$ , but  $T(t)$  retains some memory of previous states and thus is not fully determined by  $X(t)$ . For example, as shown in Figure 2, suppose  $Z_B(t)$  first rises above  $M$  on day  $t = 16$  and dips back below  $M$  on day  $t = 18$ . Then,  $T(t) = 0$  for all days through  $t = 15$ ,  $T(16) = 1$ ,  $T(17) = 2$ , and  $T(18) = 0$ . With this notation, we can say that “ $X(t)$  is the beginning of a blocking event” if

$$T(t) = 1 \quad \text{and} \quad T(t + D - 1) = D. \quad (2)$$



**Figure 1.** (a) Blocking fraction (the percent of days with  $T \geq 5$  days) for MM model data with  $M = 100$  m. (b) Total blocking event counts for MM model data during the simulation. (c) Blocking fraction for ERA5 reanalysis data with  $M = 150$ . (d) Total blocking event for ERA5 reanalysis data with  $M = 150$  m. In all subfigures, the region we focus on is indicated by the white rectangle centered at  $0^\circ\text{E}$  and  $62^\circ\text{N}$  (approximately spanned by 3 longitude points covering  $4^\circ\text{W}$ – $4^\circ\text{E}$ , and 2 latitude points covering  $60^\circ\text{N}$ – $64^\circ\text{N}$ ).



**Figure 2.** Left: The blocking persistence problem: given a nascent blocked state, the goal is to forecast whether it will persist into a long-lasting blocking event, or quickly return to climatology. The percentile represents the climatological probability. Right: A sample trajectory of  $Z_B(t)$ , the anomaly of geopotential height defined in Section 2. The vertical dashed lines indicate new blocked states ( $T = 1$ ). The red shading indicates the duration of the block. The label  $Y = 1$  indicates that the blocked state persisted 5 days to constitute a blocking event, while  $Y = 0$  indicates that it did not.

The condition  $T(t + D - 1) = D$  only holds when there are at least  $D$  consecutive days with  $Z_B(t) \geq M$  starting from  $t$ . We can see an example of this in Figure 2 at day 24, for both a block of duration 5 and 7 days. Here,  $T(24) = 1$ , and  $T(28) = 5$ , triggering the condition for  $D = 5$ . The flow remains blocked through  $T(30) = 7$ , such that day 24 would also count as the onset of a  $D = 7$  day blocking event.

With this formulation, our central question becomes: given a  $T(t) = 1$  state at time  $t$  (the flow has just become blocked), will it stay blocked for  $D$  days,  $T(t + D - 1) = D$ , or not? We address this question by estimating the conditional probability:

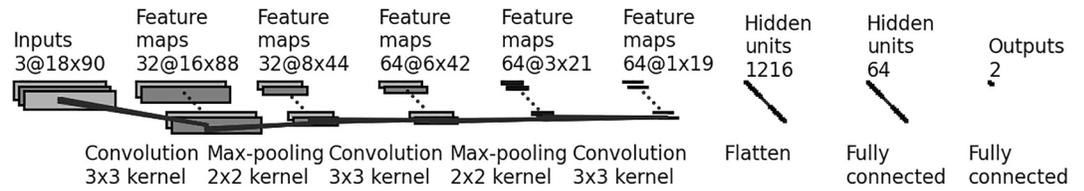
$$q(\mathbf{x}(t)) = \mathbb{P}[T(t + D - 1) = D \mid \mathbf{X}(t) = \mathbf{x}(t), T(t) = 1]. \quad (3)$$

Many recent studies have summarized extreme climate and weather events via similar functions of state (e.g., Finkel et al., 2021; Jacques-Dumas et al., 2023; Lucente et al., 2022; Miloshevich et al., 2023; Tantet et al., 2015). Unless otherwise specified, we adopt  $D = 5$  to maintain consistency with the common blocking indices (Dole & Gordon, 1983; Pelly & Hoskins, 2003; Tibaldi & Molteni, 1990). We also consider more extreme events with  $D = 7$  and  $D = 9$ .

#### 4. Convolutional Neural Network Training and Performance

Convolutional Neural Networks (CNN) have gained widespread application in probabilistic forecasting problems (Ham et al., 2019; Liu et al., 2016; Miloshevich et al., 2023) for their outstanding performance on multidimensional data sets with spatial structure. A CNN differs from a dense neural network in the use of convolutional layers with shared weights and biases across layers within the network, designed to extract features that exhibit translation invariance across the input space (Goodfellow et al., 2016). Originally developed in the context of image processing, CNN excels in scenarios where target objects, such as the face of a cat, may appear at different places within the training image. Convolutional layers allow the network to efficiently learn predictive features, combining information across multiple images. In our context, we expect predictive contributions from atmospheric eddies and Rossby waves, which share similar dynamics across all longitudes. A CNN can potentially extract these features more effectively than a fully connected architecture could, while still learning how they vary with longitude due to topography and other zonal asymmetries.

The structure of the CNN in this investigation follows Miloshevich et al. (2023) and is shown in Figure 3. It consists of a three-layer architecture, combining convolutional filters followed by ReLU activations. Specifically, we use 32 and 64 filters ( $3 \times 3$ ) for the first and last two convolutional layers. Between each pair of convolutional layers is a max-pooling layer. The output is then flattened and passed to a dense layer with 64



**Figure 3.** Convolutional Neural Network structure. The three convolutional layers respectively use 32, 64, and 64 filters ( $3 \times 3$ ), followed by ReLu activations. Between each pair of convolutional layers is a max-pooling layer with window size  $2 \times 2$ . Then the output is flattened and passed to a dense layer with 64 neurons that produces 2 outputs. A softmax function maps these outputs to two positive numbers between zero and one, representing the estimated probabilities of the nascent blocked state to persist or decay.

neurons that produces 2 outputs. Finally, a softmax function converts these two outputs to complementary probabilities.

We performed experiments with alternative CNN structures and found that reducing the widths of layers mitigates overfitting, but also reduces the performance at the best epoch (not shown). Therefore we adopt the architecture in Figure 3 and use early stopping to avoid overfitting, as detailed below.

#### 4.1. Training and Test Data Sets

We create a training and test set of all states where the flow has just become blocked:  $\{(X, T) | T = 1\}$ , where  $X$  are  $18 \times 90 \times 3$  (latitudes  $\times$  longitudes  $\times$  pressure at levels of 200 hPa, 500 hPa, 800 hPa) grid maps of geopotential height from  $20^\circ\text{N}$  to  $87^\circ\text{N}$ . Our goal is to classify which of these cases persist into blocking events ( $Y = 1$ ) versus states that do not ( $Y = 0$ ). Figure 2 shows a sample time series with 4 instances of a nascent blocked state,  $t = 16, 24, 38$  and  $47$ , only the second of which evolves into a persistent blocking event:  $Y = 0, 1, 0$ , and  $0$ , respectively. For each case, the model must classify  $Y = 0$  or  $Y = 1$  given only  $X$  at the onset time.

We examined the sensitivity of CNN model performance with respect to different amounts of training data. To prepare the data set, we integrate the MM model for 1,250k days in total. The computational cost is low, requiring 1 CPU core and approximately 11 hr. We select the first  $n$  days (with  $n$  ranging from 1k to 1,000k) to create the training data set, and always take the last 250k days for the test data set. Thus all models can be fairly compared. The trajectory length and the corresponding number of nascent blocked state states are shown in Table 1. The likelihood  $q$  of forming a blocking event varies depending on different persistence thresholds  $D$ . This dependence relationship is illustrated in Table 2.

#### 4.2. Learning Procedure

For simplicity, we use binary cross entropy as a loss function, a common choice for classification (Miloshevich et al., 2023). Alternative loss functions have been studied by Rudy and Sapsis (2023). The loss function  $L(q)$  is defined as follows:

$$L(q) = -\frac{1}{N} \sum_{i=1}^N [Y_i \log q(Y_i = 1) + (1 - Y_i) \log(1 - q(Y_i = 1))]$$

where  $q(Y_i = 1) \in (0, 1)$  is the probability of the event  $Y_i = 1$  as predicted by the CNN.  $L(q)$  is small when the CNN assigns high probability to positive events and low probability to negative events.

Given the rarity of blocking events, the data exhibit a pronounced class-imbalance, which becomes increasingly severe for longer block durations. As shown in Table 2, for  $D = 5$ , only about 1 in 5 nascent blocked states persist into an event, but  $D = 9$ , less than 1 in 20 evolve into persistent events. With this extreme imbalance, a model that never predicts an event will be correct over 80% or 95% of the time, respectively. However, such a model would clearly underperform in terms of precision and recall (defined in the next subsection), which would both be zero.

To address the class imbalance, for our results in this section we employ over-sampling (Johnson & Khoshgoftaar, 2019) techniques during training. In each epoch, we sample an equal number of nascent blocks from both

classes until we complete an iteration over all the nascent blocks in the overrepresented class. As a result, the nascent blocks that persist have been sampled multiple times during each epoch.

### 4.3. Performance Metrics

Throughout this study, we evaluate model performance using two key metrics: *precision* and *recall*. We monitor the values of these metrics on the test data set throughout the training process to determine the stopping point in order to avoid overfitting. The precision and recall are respectively defined as

$$\text{Precision} = \frac{\text{True positives}}{\text{True positives} + \text{False positives}}, \quad (4)$$

$$\text{Recall} = \frac{\text{True positives}}{\text{True positives} + \text{False negatives}}, \quad (5)$$

where “True positives” is the number of data points with  $Y = 1$  for which our CNN correctly predicts a persistent blocking event; “False positives” is the number of data points with  $Y = 0$  for which our CNN incorrectly predicts a persistent blocking event; and “False negatives” is the number of data points with  $Y = 1$  for which our CNN incorrectly predicts that the blocked state does not persist.

More informally, the precision measures the fraction of *forecasted* persistent blocks that *actually* persist. The recall, on the other hand, is the fraction of *actually* persistent blocks that are successfully *forecasted*. If one randomly predicts events with the climatological mean rate  $p$ , regardless of the system state, then the precision and recall are both given by  $\frac{p^2N}{p^2N + (1-p)pN} = p$ . This sets the floor for a useful predictor: both the precision and recall must be higher than the climatological rate.

There can be tradeoffs between improving the precision and recall. Predicting the event all the time will give you a perfect recall, but climatological precision  $p$ . A low recall implies missing a substantial number of positive events, leading to inadequate preparation and increased risk of damage. Conversely, a low precision suggests over-predicting events, “crying wolf” too often. In the context of extreme weather forecasting, this can lead to over-preparation, consequently reducing the efficiency of regular societal operations, as well as trust.

A reasonably high value of both recall and precision is crucial for an effective and resource-efficient forecasting model. We use a simplistic definition of “best” performance, expressed as

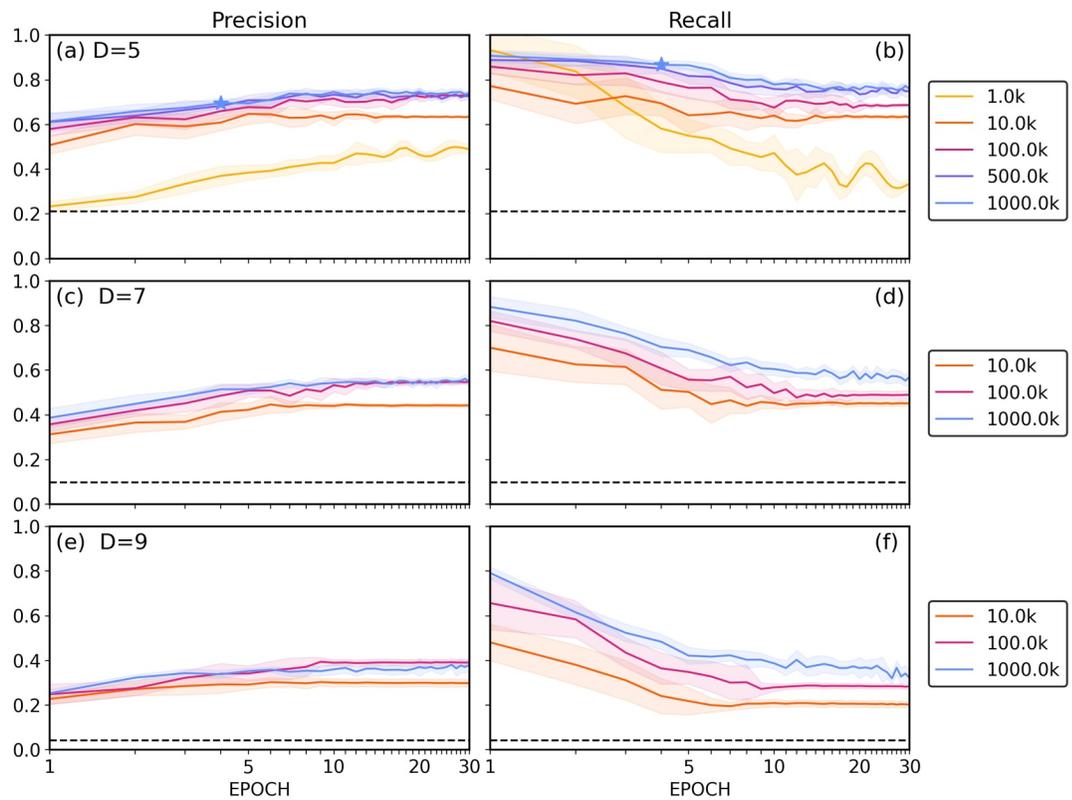
$$\text{Overall performance} = \text{Precision} + \text{Recall}. \quad (6)$$

However, it is crucial to note that in practical scenarios, designing overall performance metrics requires careful consideration of the cost of preparing versus risk of damage associated without preparation.

This naïve criterion only works when the precision and recall are both reasonably high, since forecasting the event all the time will yield a performance score of  $1 + p$  (recall of 1 and precision of  $p$ ). We used caution in ERA5 based forecasts, requiring our trained models exhibit nontrivial precision above the climatological rate. We found that the F1-score (Sasaki, 2007), another common performance metric, selects the same epoch as the metric in Equation 6.

### 4.4. Performance and Early Stopping Technique

The top row of Figure 4 shows the precision and recall evaluated on the test data for varying training data sets for  $D = 5$ . Both the precision and recall metrics are plotted starting from the end of Epoch 1 (the leftmost point on the horizontal axis of Figure 4); From Epoch 2 to Epoch 10, the precision increases, chiefly reflecting a decrease in the false positive rate, as the CNN becomes better at discriminating between persistent and non-persistent flow configurations. At the same time, the recall slowly decays: the false negative rate rises slightly as the network becomes more conservative and less likely to over-predicting persistent cases. Except for the low data regime (1k days), the performance of the CNN asymptotes after approximately 10 epochs where the precision and recall are approximately equal, but this is not necessarily the ideal stopping time (Miloshevich et al., 2023).



**Figure 4.** Precision (a, c, e) and recall (b, d, f) as a function of training epoch, for convolutional neural networks (CNNs) trained on data sets of varying sizes (curve color) and thresholds of blocking persistence (rows). As detailed in the text, all the models are tested on events from the same 250K data set not seen in training. Panels (a, b) show results for 5 days blocks ( $D = 5$ ); for example, the light blue curves are trained on all events in the full 1,000k-day simulation, while the other curves show results based on smaller training sets as indicated by the legend. The blue stars indicate the “best” convolutional neural network (see text), with a precision = 0.70 and recall = 0.87. Panels (c, d) show results for  $D = 7$  blocks and (e, f) for  $D = 9$  blocks. Fewer curves are displayed for  $D = 7$  and  $D = 9$  for the sake of clarity. Shading indicates uncertainty, assessed by taking one standard-deviation of results of 10 neural network training with i.i.d random parameter initialization.

To select the CNN parameters with the best performance, we assessed the overall performance defined in Equation 6 at the end of each epoch. We then use the parameters from the epoch with the largest value. The “best” CNN is obtained by training on the full data set of 1,000k days for 4 epochs, indicated by the star in Figure 4. It achieves precision of 0.70 and recall of 0.87, exhibiting significant predictive power over the climatological mean prediction (the black dashed line with value 0.21). Therefore, we use it for further analysis in Section 5.

All of our CNNs significantly outperformed the climatological mean prediction for any amount of data or training length. Interestingly, although the best performance is always realized with the longest trajectory of 1,000k days, precision and recall have different sensitivities to the training data size. For  $D = 5$  events, the precision improves with more data up to 100k days (equivalent to approximately 1,000 winters), after which additional data does not lead to much improvement. The recall, however, is more data-hungry; its performance continues to improve with more data up to 500k days, equivalent to 5 millennia of winter data. This reflects the fact that more data continues to help the CNN avoid missing events after its ability to limit false positive forecasts has saturated.

Figure 4 also shows the results for higher persistence thresholds,  $D = 7$  and 9. These thresholds correspond to rarer events, and even with the longest trajectory of 1,000k days, the precision and recall curves suffer for two reasons. First, as seen from Table 2, the number of positive events drops, effectively limiting the data set almost by a factor of 5 for the most extreme  $D = 9$  cases. More importantly, however, it simply becomes harder to discriminate rare events as the data set becomes more imbalanced: less than 1 in 10 nascent blocking states will evolve into a 7 days blocking event, and less than 1 in 20 into a 9 days blocking event. Without our efforts to overcome this imbalance, a network can classify almost all events correctly by never predicting a persistent case.

Despite the difficulties, the CNNs still show some skill in rare event forecasting. Given the full 1,000k data set, for  $D = 9$  the precision and recall converge to about 0.35, a factor of two worse than the CNN in the  $D = 5$  case but a factor of 10 better than climatology. As with the  $D = 5$  cases, we found that the recall for  $D = 7$  and 9 suffers more than the precision when the data set shrinks: with less events to learn from, the CNNs become more conservative and less likely to call an event. The recall depends on the false negative rate, and thus appears more sensitive to class imbalance. More data gives the network more true positive cases to learn from, apparently helping to overcome this challenge.

The low precision and recall values for smaller data sets (1k and 10k) do not bode well for training our CNN on ERA5 data, which will be discussed in detail in Section 7. For  $D = 5$ , there are 273 nascent blocked states in the ERA5 record, 84 of which persist into blocking events (see Table 3). This data amount falls between our 1k and 10k cases where data clearly limit performance. Consistent with our experience with the MM model, recall is the metric that suffers most from limited data, and stands to benefit the most from TL.

## 5. Feature Analysis: What Is Our CNN Using to Predict Blocking Events?

Before turning to forecasting in the realistic data regime, we ask what our best CNNs have learned to make these forecasts. Explainable Artificial Intelligence (XAI) is an array of techniques used to try to gain some understanding of the basis on which neural networks make predictions (Linaratos et al., 2020). In this section, we use SHapley Additive exPlanation (SHAP) value analysis to dissect the contributions of different atmospheric pressure levels and geographic areas that our CNN is using to make its predictions. We further construct a sparse model using the identified important features as inputs to quantitatively justify their relative importance in the prediction process.

### 5.1. XAI Method

SHapley Additive exPlanation (SHAP) values, introduced by Lundberg and Lee (2017) and Shrikumar et al. (2017), draw inspiration from Shapley values in game theory (Lipovetsky & Conklin, 2001). In the domain of weather and climate science, SHAP values have found broad use, with applications ranging from Earth System model error characterization (Silva et al., 2022) to drought forecasting (Dikshit & Pradhan, 2021).

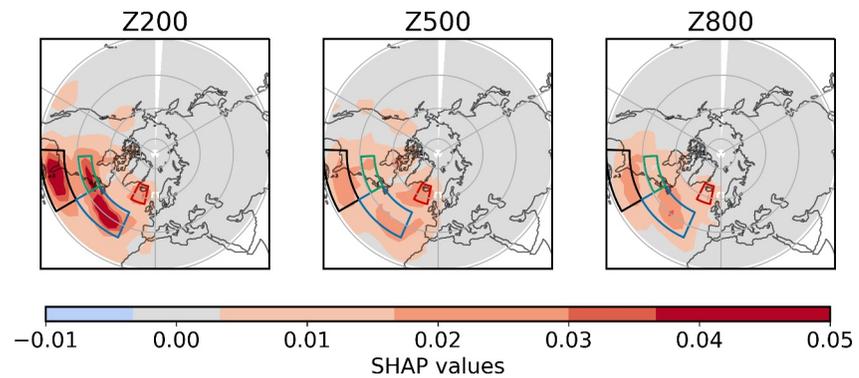
Intuitively, given a function  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  (such as the conditional probability  $q$  in Equation 3), SHAP assigns an importance value  $\phi_i$  to each feature  $x_i$  of the argument  $\mathbf{x} \in \mathbb{R}^d$ , which combine additively:

$$f(\mathbf{x}) = \mathbb{E}[f(\mathbf{x})] + \sum_{i=1}^d \phi_i(f, \mathbf{x}). \quad (7)$$

With no knowledge of  $\mathbf{x}$ , the optimal prediction of  $f$  (in a mean square sense) is the climatological average over the distribution of  $\mathbf{x}$ :  $\mathbb{E}[f(\mathbf{x})]$ . SHAP values quantify how much is gained beyond this baseline by incorporating information from each component  $i$  of  $\mathbf{x}$ . The SHAP values  $\phi_i(f, \mathbf{x})$  are unique for each sample of  $\mathbf{x}$ , but features  $i$  for which  $|\phi_i(f, \mathbf{x})|$  are large for most  $\mathbf{x}$  (i.e., a large average SHAP value) can be singled out as important, or useful, for the prediction of  $f(\mathbf{x})$ . SHAP values possess advantageous theoretical properties as well, and we refer the reader to Lundberg and Lee (2017) for a detailed theoretical analysis. In this study, SHAP values are computed using the Python package DeepSHAP (Chen, 2022). The function  $f(\mathbf{x})$  is taken as the estimated conditional probability  $\hat{q}(\mathbf{x})$  computed by the CNN, that is, the probability, according to the CNN, that the blocked state will extend  $\geq D$  days, leading to a blocking event.

### 5.2. Results

Figure 5 shows the composite of SHAP values for true positive data. Because few nascent blocks persist for  $D = 5, 7, \text{ or } 9$ , the climatological probability of a persistent event  $\mathbb{E}[\hat{q}(\mathbf{x})] = 0.21, 0.096, \text{ and } 0.044$ , respectively. For our CNN to call a positive event, we require the conditional forecast probability  $\hat{q}(\mathbf{x})$  to be larger than 0.5. Hence a positive (negative) value of  $\phi_i(\hat{q}, \mathbf{x})$  indicates that knowing the geopotential height anomaly at this level and location increases (decreases) the likelihood of a positive event. Therefore, the shading in Figure 5 can be interpreted as the average influence of each grid point for the CNN to successfully predict a long-lasting blocking



**Figure 5.** Composite maps of SHAP values,  $\bar{\phi}$ , of geopotential height at 200, 500, and 800 hPa, for true positive cases, that is, when the convolutional neural network (CNN) accurately forecasts a persistent blocking event. The unit is the probability per feature ( $Z$  at a given location and pressure level) of a positive forecast (see Equation 7), indicating the feature's average incremental contribution to the CNN's confidence that the nascent blocked state will evolve into a persistent blocking. The boundaries of the most important regions learned by the CNN are marked by solid lines and denoted region 1 (Florida, black), region 2 (north Atlantic, blue), region 3 (northeastern North America, green) and region 4 (Iceland, red).

event. For the averages over each region, the standard deviations for Z200, Z500, and Z800 are 0.039, 0.026, and 0.028, respectively, with a roughly symmetric distribution, indicating that the SHAP value analysis in Figure 5 represents the overall sample behavior, rather than being skewed by outliers.

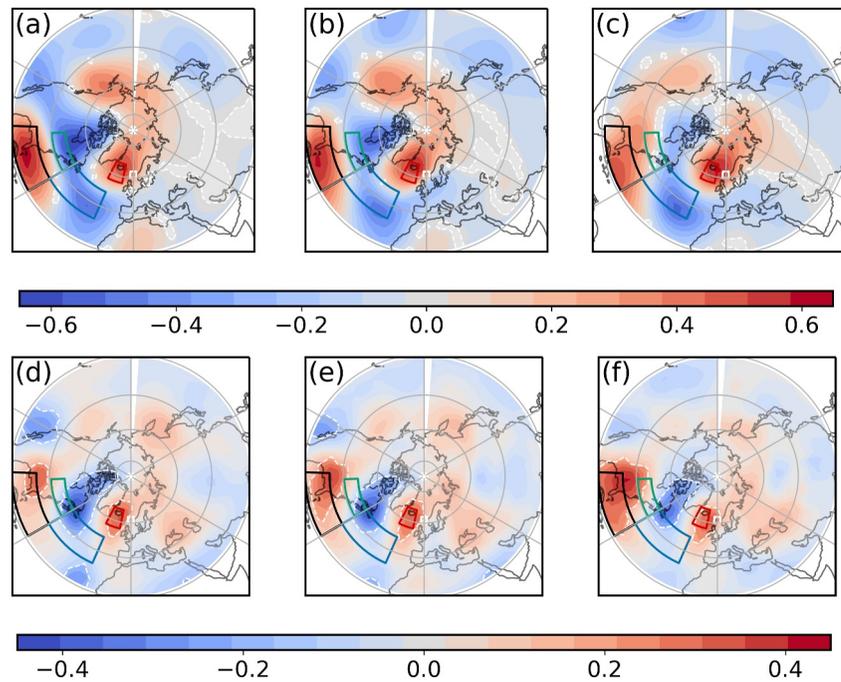
The SHAP composite is approximately uniformly non-negative because it is based only on true positive events: additional information should always increase the forecast probability. This indicates that the CNN has been well-trained to only use geopotential height information that improves the blocking event probability, and suggests it has identified robust features that herald a persistent block. A composite based on true negative cases (not shown) reveals similar patterns, but of the opposite sign.

The first thing to notice is that anomalies upstream from the blocking region (to the west) are more valuable than other regions for predicting the persistence of the blocked state. Moreover, the commonality among different pressure levels reflects the relatively barotropic nature of the MM model. In general, however, the CNN prediction relies most on the upper level flow (200 hPa).

The SHAP values emphasize four distinct regions in a quadrupole arrangement to the west of the Atlantic blocking region, as marked in Figure 5. We chose these regions to encapsulate high SHAP values using the following algorithm: after objectively identifying regions where SHAP values exceeded a set threshold, we defined boundaries by hand with the goal of enclosing these regions across all three levels within the smallest encompassing rectangle. While part of the goal of choosing these regions was to build a sparse predictor in the next section, they give us physical insight on their own.

The meaning of the SHAP values can be more easily interpreted with the aid of composites of the 3,341 true positive events (Figure 6), which show us the sign of anomalies that favor persistence. Positive geopotential anomalies in region 1 (black, centered over Florida) and 4 (red, over Iceland, just east of the blocking region itself) at the onset of blocking indicate to the CNN that a block will persist, while negative anomalies over Regions 2 (blue, North Atlantic Ocean) and 3 (green, northeast US) also favor persistence.

Regions 2 and 4 project onto opposing centers of action of the North Atlantic Oscillation (NAO). They indicate that a more negative NAO state at the onset of blocking increases the likelihood of a persistent block. Previous studies have also found that blocks tend to be more persistent when the NAO is negative (Barnes & Hartmann, 2010). While a blocking pattern off Europe projects weakly onto the NAO itself, SHAP analysis indicates that the wider structure of the pattern is important. Regions 1, 3, and 4, on the other hand, appear to be part of a wave train arching southwest from the blocking region. Their importance suggests that downstream development of a wave packet propagating along the jet stream helps drive persistent blocking events in the North Atlantic.



**Figure 6.** Averages of nascent blocking states that evolve into persistent blocking events ( $T = 1, y = 1$ ) of (top row, a–c) MM data set and (bottom row, d–f) ERA5. The colorbar represents values of geopotential height anomalies normalized by the standard deviation at each location and height. The white dashed line indicates the 0.05 significance level for a one sample  $t$ -test of the null hypothesis that the expected value is zero. The box areas identified by SHAP analysis lie in statistically significant regions. Regions that are not significant are shaded by white. For MM model data set (the top row), most of the regions are statistically significant, while for ERA5 data set (the bottom row) most of the regions are not statistically significant.

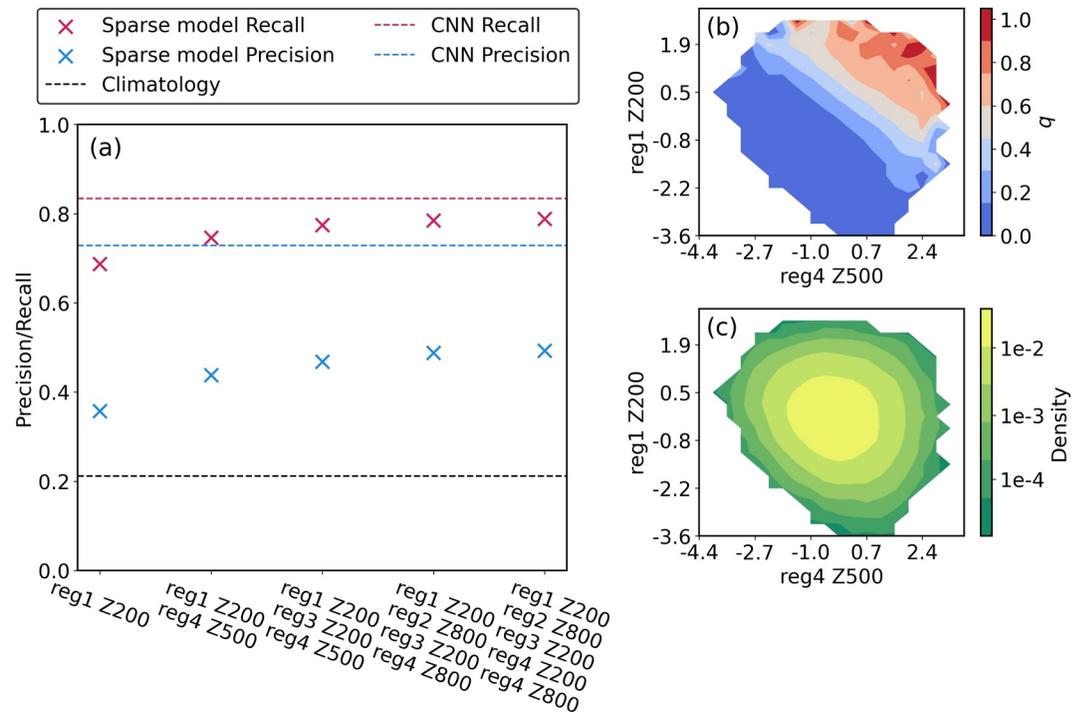
## 6. Building a Sparse Model With Logistic Regression

In quantifying the relative importance of the geopotential height as a function of location, the SHAP values suggest there is potential for dimension reduction. The CNN did not rely significantly on the information about  $Z$  in the large gray regions downstream of our target blocking area in Figure 5 to predict the potential persistence of a nascent blocking anomaly. Intuitively, conditions over central Asia will take some time to affect the flow over the North Atlantic, and are mostly irrelevant for a forecast at 5 days range.

To gain more physical insight into the utility of the SHAP values, and so gain confidence in our CNNs, we explored a simplistic dimension reduction approach focused on the regions highlighted in Figure 5. Our aim was not to achieve the ideal dimension reduction, but to provide physical insight. Thus we ask: how well can one predict the persistence of a nascent block given only the very coarse information about the flow provided by the average geopotential height within these regions at the three levels?

For these simple models, we computed the local mean of  $Z_{200}$ ,  $Z_{500}$ , and  $Z_{800}$  for each of the four rectangles shown in Figure 5, resulting in 12 time series. We then applied logistic regression with all possible combinations of these 12 features for subsets of dimension up to 5, that is, for dimension 1, fitting a logistic function with each time series alone, for dimension 2, all possible combinations of two time series, and so forth. The results for the sparse models with the best predictive skill on the test set are illustrated in Figure 7a. The horizontal axis denotes the combinations of variables that achieve the predictive scores shown in the figure.

We draw three key conclusions from Figure 7a. First, to predict the persistence of a blocked state, the best one-dimensional feature is  $Z_{200}$  in region 1, over Florida and the Gulf and upstream of the block, not  $Z_{500}$  in region 4, the  $Z$ -field nearest to the block. Second, the combination of  $Z_{200}$  in region 1 with  $Z_{500}$  in region 4 forms a two-dimension model (shown in Figure 7b) that already recovers a recall value of 0.75—it captures three quarters of all blocking events—with a precision of 0.44, twice the climatological rate. The precision and recall of the full



**Figure 7.** (a): Sparse model predictive skill on the test data set. The horizontal axis represents the dimension  $d$  of the sparse model from 1 to 5, with labels showing the combination of variables (“R1” = “region 1”) that achieves the best predictive skill among all combinations of  $d$  variables. (b) Conditional probability of a persistent block,  $q$ , as a function of mean normalized geopotential height anomaly at 200 mb over region 1 and at 500 mb over region 4 (the second column of (a)). (c) The marginal density (likelihood of observing these anomalies) as a function of the same variables. Densities below  $10^{-5}$  are cut off.

CNN, however, are 0.87 and 0.70. This leads us to the third key message: there is a large discrepancy in precision between CNN and logistic regression. Even with 5 predictors, the precision of our sparse model is only 0.5.

The poor precision indicates that the sparse model makes too many false positive predictions. This could suggest that the decay of the Atlantic blocked state is a more nonlinear dynamical phenomenon, which cannot be modeled as a simple linear statistical model. A CNN can capture these nonlinearities more effectively than sparse regression, which is consistent with previous research which found North Atlantic blocks are associated with nonlinear processes (Evans & Black, 2003). It could also indicate that more subtle features outside these 4 centers (and variation within these regions) are important. Figure 5 indicates that the CNN uses information across all of the North Atlantic, eastern North America, and even off the west coast of the US, to make skillful predictions.

To explore the effectiveness of the two-dimensional sparse model, we visualized the conditional probability of a block persisting,  $q$ , projected onto this simple subspace (shown in Figure 7b). For example, the lightest pink region, corresponding to  $q \approx 0.5$  indicates that if, at the onset of blocking,  $Z$  at 200 hPa over region 1 (Florida) is particularly high or  $Z$  at 500 hPa in region 4 (Iceland) is abnormally high, the system has a roughly 50% chance of evolving into a persistent block, more than double the climatological rate of 21%. In the red region at the top right, where both of these regions exhibit abnormally high pressure, the chance of a persistent block increases to near 100%.

Figure 7c shows the likelihood of observing these Z200 and 500 anomalies. Most often, the system exists in the middle of the diagram, where the probability of a blocking event hovers near or below the climatological value. The most likely state that exhibits a high chance of a block lies along the diagonal from the upper left to the lower right, with moderately high Z200 and 500 anomalies. The states in the top right corner, for which a persistent block is nearly certain, are very rare.

The sparse models suggest physical links between blocking events and the upstream flow. The Atlantic blocking region lies at the end of the Atlantic storm track (Michelangeli & Vautard, 1998). Persistent blocks, at least in the

MM model, are favored when there is enhanced wind off the east coast of the US (high pressure over Florida, region 1) and low pressure over regions 2 and 3 (which are highlighted in the higher dimensional sparse models). This displaces the climatological winds upstream of the blocking region equatorward. This will modify the input of storm activity into the blocking region, consistent with prior studies that have highlighted the relation between the storm track and blocking events (Yang et al., 2021; Zappa et al., 2014b).

## 7. Extending to ERA5 Using Transfer Learning

Given sufficient data, it was possible to construct a CNN that skillfully forecasts the maintenance of blocking events in the MM model. However, the ERA5 data from December, January and February (DJF) between 1940 and 2022 exhibit only 273 nascent blocked states in our Atlantic region of focus. Unfortunately, this low-data regime is where we see a significant degradation in performance in Figure 4. The curve associated with the trajectory of 10k days (699 nascent blocked states) plateaued at lower values for both the precision and recall. With only 1k days (63 nascent blocked states) performance was poor, and the learning unstable, oscillating significantly across epochs.

The class imbalance between  $Y = 0$  and  $Y = 1$  adds to the difficulty (see Table 3), particularly when longer blocks are considered. An extreme example is the set of blocking events that last  $\geq 9$  days: there are only 18 such events in the reanalysis record out of 273 data points. Such a small sample of positive data can hardly support any meaningful training, and makes it impractical to get meaningful uncertainty bounds on performance. In a standard training-test data split with a ratio of 90:10, only around 2 positive events typically fall in the test set, making it challenging to robustly assess the skill.

When training on the limited number of events in the reanalysis, a CNN can more easily suffer from overfitting, where the network uses “noise” (unrelated features) to classify blocking events. Overfitting can be diagnosed when the performance on the test set diverges from the training set. Yang and Gerber (submitted) found that the oversampling strategy used so far in this study was more prone to overfitting than a weighted loss function strategy (Johnson & Khoshgoftaar, 2019). With this latter strategy, one emphasizes the rare class (in our case, positive events) by increasing its weight in the loss function. In our remaining experiments, we weighted positive and negative events inversely to their occurrence rate.

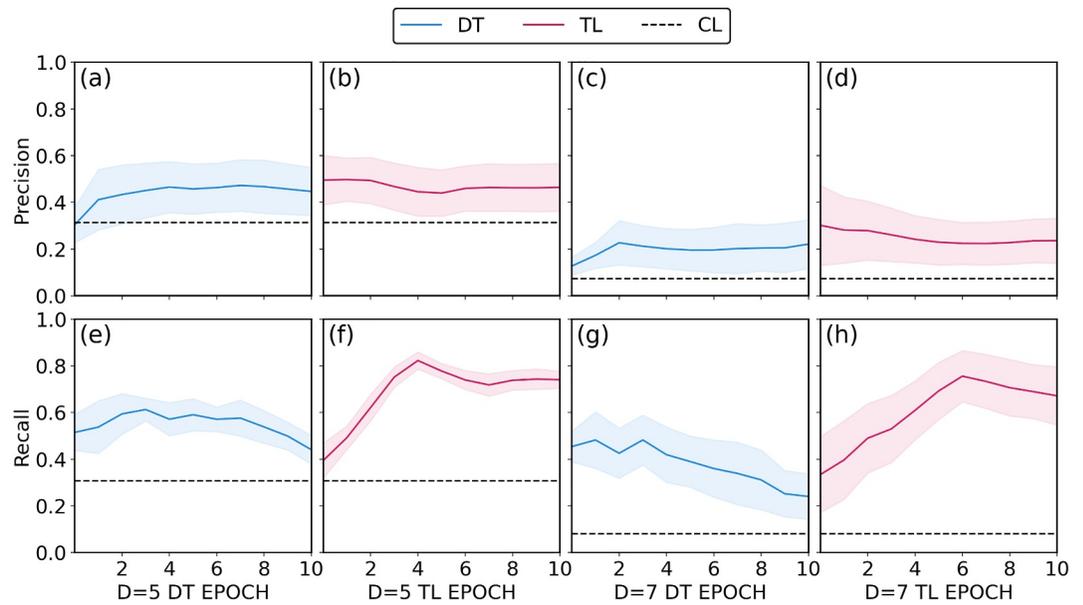
### 7.1. Direct Training

The scarcity of events makes direct training (DT) on ERA5 blocks challenging. In our study of the MM model data, we had the luxury of a large test data set (which we intentionally kept the same for fair comparison of the different CNNs), even for the case with only 1k training days. For ERA5 data, we use cross validation (Goodfellow et al., 2016) to make the best use of the smaller data set. The limited number of states were partitioned into training and test sets in ratios of 90:10; we also tried 80:20, and the results were similar (not shown). These splits were chosen to balance two difficulties: a small training set can prevent robust learning, while a small test data set limits accurate evaluation, even for a well-trained model.

To proceed, we first reduced the resolution of the ERA5 data to a comparable size of the MM output, considering geopotential height on the same three levels at the same coarse resolution. Reducing the resolution allowed us to use the same CNN architecture, and made TL possible (as discussed below). It also helped avoid overfitting, reducing the number of input variables relative to the number of events. Then we created the test-train splits, yielding 10 cross validation sets with distinct test events. Finally, for each test-train split, we trained and evaluated 10 CNNs, where variations were confined to random weight initialization and shuffling of training data.

Providing meaningful uncertainty on the precision and recall statistics from DT, shown in the left column of Figure 8, is challenging. As the 10 CNNs trained on each train-test split are not independent and identically distributed (IID), we first average the skill scores within each split. The 10 test sets, however, can be viewed as IID samples. The solid lines and shades respectively represent the mean and two-standard deviation bounds of the precision and recall, as a function of epoch, across the 10 splits.

For 5 days blocks, a CNN trained by DT can beat the climatological forecast, albeit only modestly. Given the small testing data set (27 nascent blocks, of which roughly eight persist into events), it is important not to put too much stock in the best possible performing network, for CNN can get lucky on a small sample size. The average performance quantifies the potential skill more reliably. On average, a CNN can achieve a precision of



**Figure 8.** Comparison of convolutional neural network forecast skill between direct training (DT, blue) and transfer learning (TL, red). Panels (a, b) compare the precision of DT training epoch and of TL fine-tuning epoch for  $D = 5$  (standard blocking events). (e, f) Compare the recall of DT training epoch and of TL fine-tuning epoch for  $D = 5$ . (c, d) Compare the same quantities as (a, b) for  $D = 7$ . (g, h) Compare the same quantities as (e, f) for  $D = 7$  (longer blocking events). The black dashed line indicates the climatological event rate  $p$ . The shading shows a two-standard deviation uncertainty bound, as detailed in the text.

approximately 0.45: when it calls a persistent blocking event, 4–5 out of 10 times it is correct, as compared to about 3 of 10 in the climatology. The recall was modestly better, the network only missing 4 of 10 actual events, while a climatological forecast would miss 7 of 10.

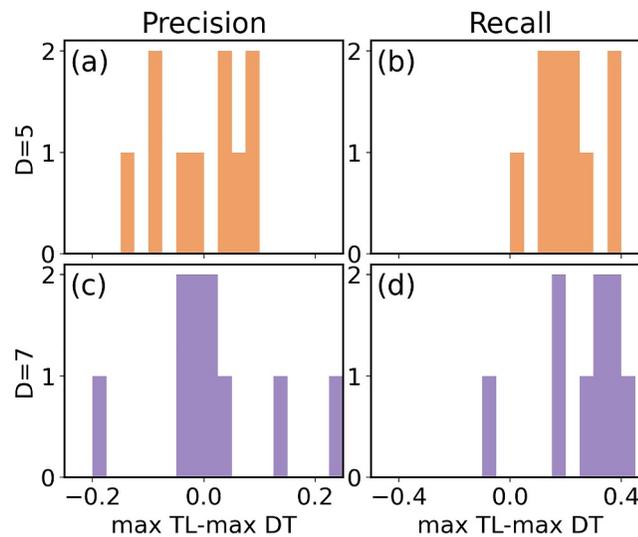
We also explore 7 days events, where only 13% of nascent blocks evolve into 7+ d events. Again, the average CNN modestly beats the climatological forecast in terms of precision: 1/5 of the cases it calls evolve into persistent events, roughly double the success rate by a guess with a Bernoulli random variable. The recall was initially deceptively high (the network captured 5 of 10 blocks), but this skill rapidly decreased with training. This was due to the fact that CNNs at early stages of DT call too many events. As it trains further, it reduces the forecast rate, declaring fewer false positives at the expense of missing more events.

## 7.2. Transfer Learning

Transfer learning has found broad application in atmospheric science, such as detecting gravity waves (González et al., 2022), improving extreme heatwave forecasts in climate models (Jacques-Dumas et al., 2022), subgrid-scale turbulence parameterization (Subel et al., 2021), image restoration (Guo et al., 2022) and parameter retrieval from raw dew point temperature profiles (Malmgren-Hansen et al., 2018).

TL involves pre-training a model on a larger data set that is similar to the data set of interest (source domain), then fine-tuning the model on the smaller target data set (target domain). This approach is particularly beneficial when labeled data for the target task is limited, as it allows the model to exploit learned features and representations from the larger data set to enhance its performance on the smaller data set. With this strength, TL has shown its power in forecasting, combining the data from a climate model (Rasp & Thuerey, 2021) or a dynamical model (Mu et al., 2020) with the observational record to improve medium-range weather forecasting and ENSO prediction.

In this section, we apply TL to leverage our MM data set to predict events in the reanalysis data. As a quasi-geostrophic model, MM has complexity between full climate models (e.g., Rasp & Thuerey, 2021) and low order models (e.g., Mu et al., 2020) used in previous TL studies. The overall process is to first “pre-train” a CNN on the MM model data set, learning to capture the characteristic features of blocking. While significantly



**Figure 9.** Histograms of the performance gap between the best performing convolutional neural networks (CNNs) obtained with Transfer learning (TL) versus the best performing CNNs obtained with direct training (DT), for precision and recall. (a) Is the performance gap of precision for 5 days events. (b) Is that of recall for 5 days events. (c) and (d) are of precision and recall for 7 days events. “Best performing” was determined by stopping the training procedure at the epoch when the best overall balance between high precision and recall was achieved in the mean (solid lines in Figure 8). The 90:10 split yields 10 different convolutional neural network scores, and the differences between pairs of TL and DT based CNNs, scored on the same test split, are shown.

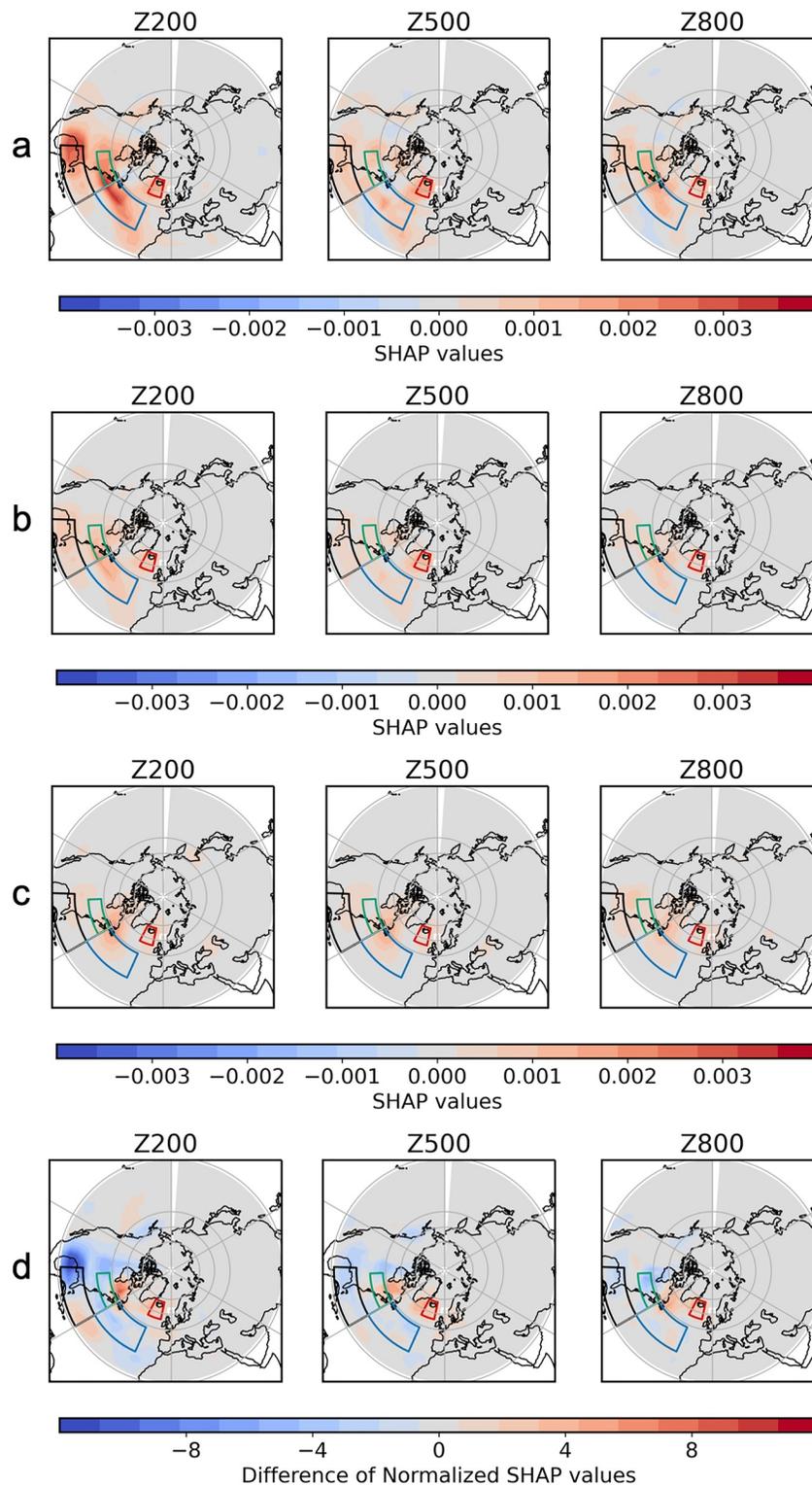
simplified, the MM model is skillful in representing atmospheric variability (Lucarini & Gritsun, 2020), but more importantly provides extensive positive and negative cases to learn from, supporting optimal CNN training, as demonstrated in Section 4. After pre-training, our CNN is then fine-tuned on the ERA5 data set, where the weights are modified to account for biases in the MM model, and the parameter scales are calibrated.

In most applications of TL, only the weights in the last few layers of a neural network are fine-tuned on the target domain (Hussain et al., 2019; Talo et al., 2019; Yosinski et al., 2014). Following this convention, we only retrain the last layer of the CNN on ERA5 while keeping the other layers frozen. This allows the CNN to correct biases it inherits from MM, but not to fall back into the poorly constrained limit we reached with DT. We also tried retraining other single layers, but retraining the last layer performed the best. To avoid overfitting, we set the learning rate to 1/10 the learning rate of pre-training.

We tested different lengths of pre-training and then evaluated the performance of the resulting models with the peak precision and recall in the transfer-learning phase. The results show that CNN parameters taken at earlier pre-training epochs show better peak performance after TL (results not shown). This suggests that overfitting on the source domain cannot be fully corrected by fine-tuning on the target domain. For the displayed results in Figures 8–10, we use a pre-training of two epochs for  $D = 5$ , and one epoch for  $D = 7$ . Given the 1,000k days of MM integration we had at our disposal, this means that the neural network has explored more than 70,000 unique nascent blocking states (all of them twice, for  $D = 5$ ) before seeing any of the 273 events in ERA5.

We follow a similar procedure as with DT to assess the ensemble-average performance. We pre-train 10 CNNs with the 1,000k-day MM data set; the only differences are due to randomness in the initialization and training data shuffling. We then carry out a 10-fold cross-validation procedure with 90:10 splits: for each split, we perform TL fine-tuning on the 10 pre-trained CNNs. We compute the mean precision and recall for each split. The results in the TL columns of Figure 8 show the mean and 2-standard deviation bounds across all the splits.

Compared to DT, TL begins with a higher precision but lower recall due to pre-training. With additional fine-tuning, the precision stays almost unchanged, while the recall grows markedly. The network is able to increase the number of events that it can capture (lowering the number of false negatives) with minimal degradation in reliability of its forecast (i.e., only slightly increasing the false positive rate).



**Figure 10.** Rows 1 through 4 are composite maps of SHAP values,  $\bar{\phi}$ , for geopotential height (200, 500, and 800 hPa), averaged over true positive predictions of blocking events in ERA5 by the convolutional neural networks (CNNs) listed below. This is the same quantity shown in Figure 5, but now applied to ERA5 events. Row *a* shows  $\bar{\phi}^{MM}$  for the pre-trained CNNs before Transfer learning (TL) (i.e., networks that have only learned from MM, but applied to ERA5). Row *b*:  $\bar{\phi}^{TL}$  of these pre-trained CNNs after fine-tuning. Row *c*:  $\bar{\phi}^{DT}$  of CNNs directly trained on ERA5 (i.e., networks that never saw the MM events). Row *d* shows the change in the SHAP values,  $\Delta\phi$ , between the first two rows, after normalization as detailed in the text. This quantifies the effect of TL: positive values indicate that information from the region became more important for the prediction, while negative values indicate that anomalies in the region became less important for prediction.

Uncertainty in the precision is dominated by differences in the true positive events between the splits; consequently, the 2-standard deviation error bounds are comparable for DT and TL. The recall is less sensitive to differences among the splits, however, and at least for the  $D = 5$  case, there is noticeably less spread across the splits with TL. This is understandable because recall, by definition, doesn't depend on the positive rate of the test data set, which varies a lot for small data sets (around 27 states in each test set after splitting). On the other hand, precision relies on the positive rate of the test data set, so it has more intrinsic variability.

We still evaluate the overall performance by Equation 6. Focusing first on  $D = 5$  events, the best mean performance with DT is a precision of 0.45 and recall of 0.61, which is realized at Epoch 3. With TL, we achieve an average performance with a similar precision of 0.45 and higher recall 0.82 (at Epoch 4). A noticeable advantage of TL is the significantly reduced variance in recall compared to DT, indicating TL's superior robustness in prediction, attributed to its enhanced capacity for capturing predictive features. For  $D = 7$  day events, the best mean performance with DT is a precision of 0.21 and recall of 0.48, achieved after three epochs. TL, however, achieves a precision of 0.22 and recall of 0.76 at Epoch 6.

To ensure that these gains in recall are statistically significant, we conducted a Wilcoxon signed-rank test (Conover, 1999). Figure 9 shows histograms of the difference in precision and recall between DT and TL. For example, each of the 10 values in the histogram for  $D = 5$  is defined for a specific train-test split, evaluated by subtracting the mean precision (recall) of 10 randomly initialized TL models taken at Epoch 4 from the mean precision (recall) of 10 randomly initialized DT models taken at Epoch 3. The spread here stems primarily from the fluctuation in 10 small-size test sets, not uncertainty in the networks due to randomness in training. The values for small-size test sets are taken at the same epoch of the best mean performance.

The average recall with TL surpasses that of DT by 34% ( $p = 0.001$ ) for 5 days events and by over 50% ( $p = 0.002$ ) for 7 days events. While there is not a significant difference between the TL and DT precision, it is critical that TL was able to improve the recall without sacrificing precision. One could easily inflate the recall by declaring more positive cases, but without any skill, the precision would suffer and approach the climatological rate.

### 7.3. What Has Transfer Learning Learned?

When we show ERA5 events to CNNs first trained on the MM data set, what exactly is the CNN learning to improve the recall? For example, do the key geographical regions and levels (Figure 5) retain the same level of significance? It is reasonable to expect that this might not be the case. In the MM data set, the duration of the Atlantic blockings could be related to upstream flow, specifically to the structure of the wave train at the blocking onset. The mechanism for blocking in the real world is more complicated, and the correlated pattern may shift, intensify, and/or weaken. To address these questions, we compare the SHAP values of the pre-trained CNNs when directly applied to ERA5 (i.e., without fine-tuning) to the SHAP values of the CNN after four epochs of fine-tuning, as shown in row *a* and row *b* of Figure 10. The most evident difference after fine-tuning is a decrease in the amplitude of the SHAP values. This is because the climatological rate of positive blocking events in ERA5 is higher: almost 1/3 of nascent blocked states persist for 5 days in ERA5, compared to about 1/5 in MM. As the expected fraction of events is larger,  $\hat{q}(\mathbf{x}) - \mathbb{E}[\hat{q}(\mathbf{x})]$  from Equation 7 will be smaller, and the SHAP value increments  $\phi_i(\hat{q}, \mathbf{x})$  will tend to be smaller. It is the sum of the SHAP values that build up the probability for a  $Y = 1$  prediction; for a more likely event, one does not need to build up the probability as much, so fine-tuning quickly adjusts the weights.

To assess the more subtle change in the relative contribution of each feature on the predicted result after TL, we show the difference in the normalized composite map  $\Delta\phi$  in row *d* of Figure 10.  $\Delta\phi$  is defined for each input  $i$  (i.e., geopotential height  $Z$  at a particular latitude, longitude, and pressure level) by 
$$\Delta\phi_i \equiv \max\left(\frac{\bar{\phi}_i^{\text{TL}}}{\sum_{j=1}^d \bar{\phi}_j^{\text{TL}}}, 0\right) - \max\left(\frac{\bar{\phi}_i}{\sum_{j=1}^d \bar{\phi}_j}, 0\right).$$
 The maximum function is used to avoid spurious negative SHAP values, which should not arise in a composite of true positive events, as discussed in the context of Figure 5. The normalization makes the total integral of the SHAP values the same for both cases, so that one can focus on where the CNN is using information, as opposed to the overall reduction of the SHAP values driven by the difference in rates.

The “normalized” SHAP values increase mainly in region 4 (the region right around the block), and additionally over Quebec and Atlantic Canada, a region less used for predictions with the MM model. The SHAP values decrease in a relative sense over regions 1 (Florida and the Gulf), 2 (North Atlantic Ocean), 3 (northeastern North America), and central North America. This change in relative importance reveals a general de-emphasis of the regions farther upstream and an increased emphasis on regions more immediately upstream. This indicates that while it is still upstream information that is most important for predicting a persistent blocking state in ERA5, the structure and westward extension of the wave train has changed.

For further insight, we compare the SHAP value patterns with a more traditional method for understanding predictability: composite analysis. Figure 6 shows composite maps of nascent blocks that evolve into persistent events in the MM model and ERA5. Persistent blocks are associated with wave activity south and west of the blocking region in both the model and reanalysis, but the pattern shifts. The wave train in MM initially arcs westward before turning southward, with a strong center of high pressure east of Florida, while the wave train in ERA5 arcs more to southwest at first, then further westward.

The SHAP values change over Quebec, capturing this shift in the wave train, but overall the CNN seems to shift to more local information with TL. We speculated that the dry, quasi-geostrophic MM model overemphasizes long-range teleconnections. It only captures deformation scale dynamics, and this only at low resolution, and so lacks smaller, local modes of instability, for example, instability associated with latent heat release due to precipitation, present in our atmosphere. The CNN makes more use of these local features when predicting the persistence of blocks, but still focuses on the upstream flow, consistent with our intuition.

Finally, we contrast the feature importance analysis of the CNN with TL (Figure 10 row *b*) to that of the CNNs trained only directly on the ERA5 output (Figure 10 row *c*). DT struggles to develop nuanced features with limited data. The SHAP values with DT are also more barotropic than those with TL. Moreover, in general, the SHAP values with TL capture finer details across a wider spatial range, while the SHAP values with DT are more localized. Geopotential height anomalies over Iceland, especially in the Z500 map, are more emphasized for TL than DT. The same applies to upstream anomalies over Florida and the Gulf of Mexico in the Z200 map. Additionally, the importance of geopotential height anomalies over the Atlantic, immediately upstream of the target region west of north Africa, is neglected in DT, though it appears in TL. This is closely correlated to the blocking event prediction from the ERA5 composite in Figure 6, which does not show as strong composite Atlantic anomaly as in the MM model.

In summary, the superiority of CNNs trained with TL, as compared to DT, appears to lie in their ability to leverage learned features from the pre-trained data set, helping the network to take advantage of information further upstream of the blocking region. In either case the precision is modest: when the networks call an event, the rate of success is at best 50% higher than a naïve climatological forecast. Pre-training the network, however, has a significant impact on the recall, increasing the forecast rate to capture more events without decreasing the precision.

## 8. Conclusion

The impact of data-driven science on weather and climate science has grown substantially in recent years. In this paper, we suggest two data-driven approaches to help predict and understand atmospheric blocking events. First, given sufficient data, CNNs are capable of identifying subtle features that differentiate short-lived blocked states from those that persist for an extended period. Moreover, Explainable Artificial Intelligence methods, like SHAP feature importance analysis, can provide insight into what features matter most to this differentiation. Second, TL has the potential to make data-driven forecasts possible for our atmosphere, making the most of the limited extreme events in the observational record by leveraging insight from longer, albeit imperfect, numerical simulations.

We began in a data-rich regime with the idealized Marshall-Molteni model, showing that a CNN can accurately predict the persistence of North Atlantic blocks in terms of both precision and recall. Leveraging SHAP feature importance analysis, we identified crucial regions for the prediction of persistent blocked states, given a nascent high-pressure anomaly. Our results suggest that incorporation of both local and non-local features is important for prediction skill.

To validate our discovery, we constructed a two-dimensional model that used only upstream anomalies over Florida and the Gulf of Mexico, and anomalies immediately upstream of the blocking region. The sparse model exhibited precision significantly above the climatological rate and recall nearly as good as the full CNN. It struggled, however, with false positives (and hence exhibited low precision relative to the CNN) which could not be improved within the log linear logistic regression framework. This suggests the CNN learns non-trivial relations in the upstream flow, extending all the way to the Pacific, to better discriminate between short-lived and long-lived blocks.

The challenge of conducting DT on ERA5 data stems from the paucity of available events. Small training and test data sets make training and evaluation difficult. With the MM model, we observed a systematic degradation in forecast skill when the training data was limited, particularly for the recall statistic. Through TL, we leverage the abundance of data generated by simplified dynamical models to enhance real-world forecasting. By pre-training a CNN on the MM model data set and retraining the deepest layer on the ERA5 data set, the recall was improved by 34% compared to a CNN developed with DT alone for 5 days events, and over 50% for more extreme 7 days events, without any loss of precision.

In addition to advancing predictive skill, TL in combination with SHAP analysis allowed us to compare the predictive features between weather systems in ERA5 and the idealized QG model. The bottom row of Figure 6 reveals biases in the MM model, which appears overly dependent on upstream features over Florida and the Gulf of Mexico relative to blocks in ERA5. This approach provides a new angle of how a machine learning approach could guide the diagnosis and quantification of model biases. This said, the success of TL results underscores the MM model's ability, despite its simplicity, to capture features that are important for predicting the persistence of blocked states in the real world. We believe that greater strides could be made by pre-training on a more advanced climate model, or even hindcasts in the subseasonal-to-seasonal (S2S) data set (Finkel et al., 2023; Vitart et al., 2017). We expect our results will help inform large-scale efforts to incorporate AI into operational forecasts, such as the AIFS model (Lang et al., 2024), which already employs TL in a different form.

The methods presented here are not limited to the context of blocking events, and can be generalized to the study of other challenging natural phenomena, especially in scenarios where data may be limited, and the potential influencing factors are complex (e.g., heat domes (Li et al., 2024)). An immediate future goal is to push further on the physical and dynamical mechanisms that causes the differences in prediction mechanisms for ERA5 and MM model. Another goal is to adapt the present approach to investigate the statistical behavior and mechanisms for the onset of the blocking events.

## Appendix A: Marshall-Molteni Model

The Marshall-Molteni (MM) model state is specified by potential vorticity  $q_j$  in three layers of the atmosphere,  $j = 1, 2, 3$ , corresponding to pressure levels 200, 500, and 800 hPa.  $q_j$  evolves according to quasi-geostrophic dynamics as

$$\partial_t q_j + J(\psi_j, q_j) = -D_j + S_j \quad (\text{A1})$$

where  $\psi_j$  is the streamfunction in layer  $j$ , related to  $q_j$  as

$$q_1 = \Delta\psi_1 - (\psi_1 - \psi_2)/R_1^2 + f \quad (\text{A2})$$

$$q_2 = \Delta\psi_2 + (\psi_1 - \psi_2)/R_1^2 - (\psi_2 - \psi_3)/R_2^2 + f \quad (\text{A3})$$

$$q_3 = \Delta\psi_3 + (\psi_2 - \psi_3)/R_2^2 + f(1 + h/H_0). \quad (\text{A4})$$

Here,  $\Delta$  is the horizontal Laplacian operator,  $R_1 = 761$  km and  $R_2 = 488$  km are the Rossby deformation radii in layers 1 and 2,  $f = 2\Omega \cos \phi$  is the latitude-dependent Coriolis parameter, and  $h$  is the orography of the surface, rescaled by the constant  $H_0$ . The operator  $D_j$  combines all dissipative terms, including radiative damping, surface friction and hyper-diffusion to crudely parametrize small scale diffusion, but is also necessary for numerical stability:

$$\begin{aligned}
 -D_1 &= (\psi_1 - \psi_2)/(\tau_R R_1^2) - R^8 \Delta^4 q_1 / (\tau_H \lambda_{\max}^4) \\
 -D_2 &= -(\psi_1 - \psi_2)/(\tau_R R_1^2) + (\psi_2 - \psi_3)/(\tau_R R_2^2) - R^8 \Delta^4 q_2' / (\tau_H \lambda_{\max}^4) \\
 -D_3 &= -(\psi_2 - \psi_3)/(\tau_R R_2^2) - EK_3 - R^8 \Delta^4 q_3' / (\tau_H \lambda_{\max}^4).
 \end{aligned} \tag{A5}$$

The forcing,  $S_j$  is computed from observed data to inject energy into the system and give the model a realistic mean state:

$$S_j = \overline{J(\psi_j, q_j)} + \overline{D_j} \tag{A6}$$

The data to construct  $S_j$  were drawn from the 1983–1992 winter (DJF) climatology of the ERA40 reanalysis provided by ECMWF.

## Appendix B: Acronyms and Definitions

Here we list the important acronyms and definitions in this paper for the convenience of the readers.

- CNN: Convolutional Neural Network - A type of deep learning model particularly effective for analyzing visual data, using convolutional layers to automatically detect and learn patterns.
- SHAP: Shapley Additive ExPlanation - A method to explain the output of machine learning models by attributing contributions of individual features based on cooperative game theory.
- MM: Marshall-Molteni - Refers to the 3-layer QG model by (Marshall & Molteni, 1993) related to atmospheric dynamics, often used in the context of studying large-scale weather patterns and teleconnections.
- QG: Quasi-Geostrophic - A simplified model in geophysical fluid dynamics that describes large-scale atmospheric and oceanic flows, assuming a balance between pressure gradient and Coriolis forces.
- XAI: Explainable Artificial Intelligence - A subfield of AI focused on making the outputs and processes of machine learning models transparent and understandable to humans.
- DG: Dole & Gordon index (Dole & Gordon, 1983)- An index developed by Dole and Gordon to quantify atmospheric blocking events, which are large-scale pressure systems that can disrupt normal weather patterns.
- DT: Direct Training - A machine learning approach where a model is trained directly on a specific data set without additional pre-training or TL techniques.
- TL: Transfer Learning - A machine learning technique where a pre-trained model is adapted to a new but related task, leveraging the knowledge gained from the original task to improve performance.
- Z: Geopotential height.
- $Z_B(t)$ : Anomalous geopotential height in our target blocking region in the North Atlantic, shown in Figure 1.
- T: Number of consecutive days of a blocked state.
- M: Threshold of geopotential height anomaly in blocking events criteria.
- D: Threshold of consecutive days in blocking events criteria.
- X: Full model state vector.
- Y: Indicator of whether a blocked state persisted.
- $q(x(t))$ : Conditional probability that a blocked state  $x(t)$  will persist.
- $L(q)$ : Binary cross entropy loss function used for classification problem.
- Precision:  $\frac{\text{True positives}}{\text{True positives} + \text{False positives}}$
- Recall:  $\frac{\text{True positives}}{\text{True positives} + \text{False negatives}}$

## Data Availability Statement

The data from the Marshall-Molteni model were generated using a Fortran code provided by Valerio Lucarini and Andrey Gritsun (Lucarini & Gritsun, 2020). The Fortran code, along with the Python code for computing SHAP values, TL and producing plots is publicly available in the open repository (Zhang, 2024). SHAP values were computed using the Python package DeepSHAP (Chen, 2022). The ERA5 reanalysis data sets from ECMWF were used for data preprocessing and ML model training and testing (Hersbach et al., 2020).

**Acknowledgments**

We thank Valerio Lucarini and Andrey Gritsun for sharing their Marshall-Molteni Fortran code. We also thank Pedram Hassanzadeh and the anonymous reviewers for many helpful comments and suggestions that strengthened the paper. This work was supported by the Army Research Office, Grant W911NF-22-2-0124. EPG acknowledges support from the National Science Foundation through award OAC-2004572. J. F. is supported through the MIT Climate Grand Challenge on Weather and Climate Extremes, and the Virtual Earth Systems Research Institute at Schmidt Sciences.

**References**

Barnes, E. A., & Hartmann, D. L. (2010). Dynamical feedbacks and the persistence of the NAO. *Journal of the Atmospheric Sciences*, 67(3), 851–865. <https://doi.org/10.1175/2009JAS3193.1>

Berckmans, J., Woollings, T., Demory, M.-E., Vidale, P.-L., & Roberts, M. (2013). Atmospheric blocking in a high resolution climate model: Influences of mean state, orography and eddy forcing. *Atmospheric Science Letters*, 14(1), 34–40. <https://doi.org/10.1002/asl2.412>

Chan, P.-W., Hassanzadeh, P., & Kuang, Z. (2019). Evaluating indices of blocking anticyclones in terms of their linear relations with surface hot extremes. *Geophysical Research Letters*, 46(9), 4904–4912. <https://doi.org/10.1029/2019GL083307>

Charney, J. G., & DeVore, J. G. (1979). Multiple flow equilibria in the atmosphere and blocking. *Journal of the Atmospheric Sciences*, 36(7), 1205–1216. [https://doi.org/10.1175/1520-0469\(1979\)036<1205:mfeita>2.0.co;2](https://doi.org/10.1175/1520-0469(1979)036<1205:mfeita>2.0.co;2)

Chen, H. (2022). suinleelab/deepshap: Nature communications code [Software]. *Zenodo*. <https://doi.org/10.5281/zenodo.6585445>

Conover, W. J. (1999). *Practical nonparametric statistics* (Vol. 350). John Wiley & Sons.

d’Andrea, F., Tibaldi, S., Blackburn, M., Boer, G., Déqué, M., Dix, M., et al. (1998). Northern Hemisphere atmospheric blocking as simulated by 15 atmospheric general circulation models in the period 1979–1988. *Climate Dynamics*, 14(6), 385–407. <https://doi.org/10.1007/s003820050230>

Davini, P., & D’Andrea, F. (2016). Northern hemisphere atmospheric blocking representation in global climate models: Twenty years of improvements? *Journal of Climate*, 29(24), 8823–8840. <https://doi.org/10.1175/jcli-d-16-0242.1>

Davini, P., & D’Andrea, F. (2020). From CMIP3 to CMIP6: Northern hemisphere atmospheric blocking simulation in present and future climate. *Journal of Climate*, 33(23), 10021–10038. <https://doi.org/10.1175/JCLI-D-19-0862.1>

Davini, P., Weisheimer, A., Balmaseda, M., Johnson, S. J., Molteni, F., Roberts, C. D., et al. (2021). The representation of winter northern hemisphere atmospheric blocking in ECMWF seasonal prediction systems. *Quarterly Journal of the Royal Meteorological Society*, 147(735), 1344–1363. <https://doi.org/10.1002/qj.3974>

Dikshit, A., & Pradhan, B. (2021). Explainable AI in drought forecasting. *Machine Learning with Applications*, 6, 100192. <https://doi.org/10.1016/j.mlwa.2021.100192>

Dole, R. M., & Gordon, N. D. (1983). Persistent anomalies of the extratropical Northern Hemisphere wintertime circulation: Geographical distribution and regional persistence characteristics. *Monthly Weather Review*, 111(8), 1567–1586. [https://doi.org/10.1175/1520-0493\(1983\)111<1567:paoten>2.0.co;2](https://doi.org/10.1175/1520-0493(1983)111<1567:paoten>2.0.co;2)

Evans, K. J., & Black, R. X. (2003). Piecewise tendency diagnosis of weather regime transitions. *Journal of the Atmospheric Sciences*, 60(16), 1941–1959. [https://doi.org/10.1175/1520-0469\(2003\)060\(1941:PTDOWR\)2.0.CO;2](https://doi.org/10.1175/1520-0469(2003)060(1941:PTDOWR)2.0.CO;2)

Ferranti, L., Corti, S., & Janousek, M. (2015). Flow-dependent verification of the ECMWF ensemble over the Euro-Atlantic sector. *Quarterly Journal of the Royal Meteorological Society*, 141(688), 916–924. <https://doi.org/10.1002/qj.2411>

Finkel, J., Webber, R. J., Gerber, E. P., Abbot, D. S., & Weare, J. (2021). Learning forecasts of rare stratospheric transitions from short simulations. *Monthly Weather Review*, 149(11), 3647–3669. <https://doi.org/10.1175/MWR-D-21-0024.1>

Finkel, J., Webber, R. J., Gerber, E. P., Abbot, D. S., & Weare, J. (2023). Data-driven transition path analysis yields a statistical understanding of sudden stratospheric warming events in an idealized model. *Journal of the Atmospheric Sciences*, 80(2), 519–534. <https://doi.org/10.1175/JAS-D-21-0213.1>

González, J. L., Chapman, T., Chen, K., Nguyen, H., Chambers, L., Mostafa, S. A., et al. (2022). Atmospheric gravity wave detection using transfer learning techniques. In *2022 IEEE/ACM International Conference on Big Data Computing, Applications and Technologies (BDCAT)* (pp. 128–137). <https://doi.org/10.1109/BDCAT56447.2022.00023>

Goodfellow, I., Bengio, Y., & Courville, A. (2016). *Deep learning*. MIT Press. Retrieved from <http://www.deeplearningbook.org>

Guo, Y., Wu, X., Qing, C., Su, C., Yang, Q., & Wang, Z. (2022). Blind restoration of images distorted by atmospheric turbulence based on deep transfer learning. *Photonics*, 9(8), 582. <https://doi.org/10.3390/photonics9080582>

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>

Hersbach, H., Bell, B., Berrisford, P., Hirahara, S., Horányi, A., Muñoz-Sabater, J., et al. (2020). The era5 global reanalysis. *Quarterly Journal of the Royal Meteorological Society*, 146(730), 1999–2049. <https://doi.org/10.1002/qj.3803>

Hoskins, B. J., James, I. N., & White, G. H. (1983). The shape, propagation and mean-flow interaction of large-scale weather systems. *Journal of the Atmospheric Sciences*, 40(7), 1595–1612. [https://doi.org/10.1175/1520-0469\(1983\)040\(1595:TSPAMF\)2.0.CO;2](https://doi.org/10.1175/1520-0469(1983)040(1595:TSPAMF)2.0.CO;2)

Hussain, M., Bird, J. J., & Faria, D. R. (2019). A study on CNN transfer learning for image classification. In *Advances in computational intelligence systems: Contributions presented at the 18th UK workshop on computational intelligence, September 5-7, 2018* (pp. 191–202).

Jacques-Dumas, V., Ragone, F., Borgnat, P., Abry, P., & Bouchet, F. (2022). Deep learning-based extreme heatwave forecast. *Frontiers in Climate*, 4. <https://doi.org/10.3389/fclim.2022.789641>

Jacques-Dumas, V., van Westen, R. M., Bouchet, F., & Dijkstra, H. A. (2023). Data-driven methods to estimate the committor function in conceptual ocean models. *Nonlinear Processes in Geophysics*, 30(2), 195–216. <https://doi.org/10.5194/npg-30-195-2023>

Johnson, J. M., & Khoshgoftaar, T. M. (2019). Survey on deep learning with class imbalance. *Journal of Big Data*, 6(1), 1–54. <https://doi.org/10.1186/s40537-019-0192-5>

Kautz, L.-A., Martius, O., Pfahl, S., Pinto, J. G., Ramos, A. M., Sousa, P. M., & Woollings, T. (2022). Atmospheric blocking and weather extremes over the Euro-Atlantic sector – A review. *Weather and Climate Dynamics*, 3(1), 305–336. <https://doi.org/10.5194/wcd-3-305-2022>

Labe, Z. M., & Barnes, E. A. (2021). Detecting climate signals using explainable AI with single-forcing large ensembles. *Journal of Advances in Modeling Earth Systems*, 13(6), e2021MS002464. <https://doi.org/10.1029/2021MS002464>

Lang, S., Alexe, M., Chantry, M., Dramsch, J., Pinault, F., Raoult, B., et al. (2024). Aifs - Ecmwf’s data-driven forecasting system. Retrieved from <https://arxiv.org/abs/2406.01465>

Li, X., Mann, M. E., Wehner, M. F., Rahmstorf, S., Petri, S., Christiansen, S., & Carrillo, J. (2024). Role of atmospheric resonance and land-atmosphere feedbacks as a precursor to the June 2021 pacific northwest heat dome event. *Proceedings of the National Academy of Sciences*, 121(4), e2315330121. <https://doi.org/10.1073/pnas.2315330121>

Linaratos, P., Papastefanopoulos, V., & Kotsiantis, S. (2020). Explainable AI: A review of machine learning interpretability methods. *Entropy*, 23(1), 18. <https://doi.org/10.3390/e23010018>

Lipovetsky, S., & Conklin, M. (2001). Analysis of regression in game theory approach. *Applied Stochastic Models in Business and Industry*, 17(4), 319–330. <https://doi.org/10.1002/asmb.446>

Liu, Y., Racah, E., Prabhat, Correa, J., Khosrowshahi, A., Lavers, D., et al. (2016). Application of deep convolutional neural networks for detecting extreme weather in climate datasets.

- Lucarini, V., & Gritsun, A. (2020). A new mathematical framework for atmospheric blocking events. *Climate Dynamics*, 54(1–2), 575–598. <https://doi.org/10.1007/s00382-019-05018-2>
- Lucente, D., Herbert, C., & Bouchet, F. (2022). Commitor functions for climate phenomena at the predictability margin: The example of El Niño Southern Oscillation in the Jin and Timmermann model. *Journal of the Atmospheric Sciences*, 79(9), 2387–2400. <https://doi.org/10.1175/JAS-D-22-0038.1>
- Lundberg, S. M., & Lee, S.-I. (2017). A unified approach to interpreting model predictions. *Advances in Neural Information Processing Systems*, 30.
- Lupo, A. R. (2021). Atmospheric blocking events: A review. *Annals of the New York Academy of Sciences*, 1504(1), 5–24. <https://doi.org/10.1111/nyas.14557>
- Lupo, A. R., Mokhov, I. I., Akperov, M. G., Chernokulsky, A. V., & Athar, H. (2012). A dynamic analysis of the role of the planetary-and synoptic-scale in the summer of 2010 blocking episodes over the European part of Russia. *Advances in Meteorology*, 2012, 1–11. <https://doi.org/10.1155/2012/584257>
- Malmgren-Hansen, D., Nielsen, A. A., Laparra, V., & Valls, G. C. (2018). Transfer learning with convolutional networks for atmospheric parameter retrieval. In *IGARSS 2018 - 2018 IEEE International Geoscience and Remote Sensing Symposium* (pp. 2111–2114). <https://doi.org/10.1109/IGARSS.2018.8518097>
- Marshall, J., & Molteni, F. (1993). Toward a dynamical understanding of planetary-scale flow regimes. *Journal of the Atmospheric Sciences*, 50(12), 1792–1818. [https://doi.org/10.1175/1520-0469\(1993\)050<1792:taduop>2.0.co;2](https://doi.org/10.1175/1520-0469(1993)050<1792:taduop>2.0.co;2)
- Matsueda, M. (2009). Blocking predictability in operational medium-range ensemble forecasts. *SOLA*, 5, 113–116. <https://doi.org/10.2151/sola.2009-029>
- McWilliams, J. C. (1980). An application of equivalent modons to atmospheric blocking. *Dynamics of Atmospheres and Oceans*, 5(1), 43–66. [https://doi.org/10.1016/0377-0265\(80\)90010-X](https://doi.org/10.1016/0377-0265(80)90010-X)
- Michelangeli, P.-A., & Vautard, R. (1998). The dynamics of Euro-Atlantic blocking onsets. *Quarterly Journal of the Royal Meteorological Society*, 124(548), 1045–1070. <https://doi.org/10.1002/qj.49712454803>
- Miloshevich, G., Cozian, B., Abry, P., Borgnat, P., & Bouchet, F. (2023). Probabilistic forecasts of extreme heatwaves using convolutional neural networks in a regime of lack of data. *Physical Review Fluids*, 8(4), 040501. <https://doi.org/10.1103/PhysRevFluids.8.040501>
- Mu, B., Ma, S., Yuan, S., & Xu, H. (2020). Applying convolutional LSTM network to predict El Niño events: Transfer learning from the data of dynamical model and observation. In *2020 IEEE 10th International Conference on Electronics Information and Emergency Communication (ICEIEC)* (pp. 215–219). <https://doi.org/10.1109/ICEIEC49280.2020.9152317>
- Mullen, S. L. (1987). Transient eddy forcing of blocking flows. *Journal of the Atmospheric Sciences*, 44(1), 3–22. [https://doi.org/10.1175/1520-0469\(1987\)044<0003:tefobf>2.0.co;2](https://doi.org/10.1175/1520-0469(1987)044<0003:tefobf>2.0.co;2)
- Pelly, J. L., & Hoskins, B. J. (2003). A new perspective on blocking. *Journal of the Atmospheric Sciences*, 60(5), 743–755. [https://doi.org/10.1175/1520-0469\(2003\)060<0743:anpob>2.0.co;2](https://doi.org/10.1175/1520-0469(2003)060<0743:anpob>2.0.co;2)
- Rampal, N., Gibson, P. B., Sood, A., Stuart, S., Fauchereau, N. C., Brandolino, C., et al. (2022). High-resolution downscaling with interpretable deep learning: Rainfall extremes over New Zealand. *Weather and Climate Extremes*, 38, 100525. <https://doi.org/10.1016/j.wace.2022.100525>
- Rasp, S., & Thuerey, N. (2021). Data-driven medium-range weather prediction with a resnet pretrained on climate simulations: A new model for weatherbench. *Journal of Advances in Modeling Earth Systems*, 13(2), e2020MS002405. <https://doi.org/10.1029/2020MS002405>
- Rex, D. F. (1950). Blocking action in the middle troposphere and its effect upon regional climate. *Tellus*, 2(3), 196–211. <https://doi.org/10.1111/j.2153-3490.1950.tb00331.x>
- Rudy, S. H., & Sapsis, T. P. (2023). Output-weighted and relative entropy loss functions for deep learning precursors of extreme events. *Physica D: Nonlinear Phenomena*, 443, 133570. <https://doi.org/10.1016/j.physd.2022.133570>
- Sasaki, Y. (2007). The truth of the F-measure. Retrieved from [https://nicolasshu.com/assets/pdf/Sasaki\\_2007\\_The%20Truth%20of%20the%20F-measure.pdf](https://nicolasshu.com/assets/pdf/Sasaki_2007_The%20Truth%20of%20the%20F-measure.pdf)
- Scaife, A. A., Woollings, T., Knight, J., Martin, G., & Hinton, T. (2010). Atmospheric blocking and mean biases in climate models. *Journal of Climate*, 23(23), 6143–6152. <https://doi.org/10.1175/2010jcli3728.1>
- Shrikumar, A., Greenside, P., & Kundaje, A. (2017). Learning important features through propagating activation differences. In *International conference on machine learning* (pp. 3145–3153).
- Shutts, G. (1983). The propagation of eddies in diffluent jetstreams: Eddy vorticity forcing of ‘blocking’ flow fields. *Quarterly Journal of the Royal Meteorological Society*, 109(462), 737–761. <https://doi.org/10.1002/qj.49710946204>
- Silva, S. J., Keller, C. A., & Hardin, J. (2022). Using an explainable machine learning approach to characterize Earth System model errors: Application of SHAP analysis to modeling lightning flash occurrence. *Journal of Advances in Modeling Earth Systems*, 14(4), e2021MS002881. <https://doi.org/10.1029/2021ms002881>
- Subel, A., Chattopadhyay, A., Guan, Y., & Hassanzadeh, P. (2021). Data-driven subgrid-scale modeling of forced Burgers turbulence using deep learning with generalization to higher Reynolds numbers via transfer learning. *Physics of Fluids*, 33(3). <https://doi.org/10.1063/5.0040286>
- Talo, M., Baloglu, U. B., Yıldırım, Ö., & Acharya, U. R. (2019). Application of deep transfer learning for automated brain abnormality classification using MR images. *Cognitive Systems Research*, 54, 176–188. <https://doi.org/10.1016/j.cogsys.2018.12.007>
- Tantet, A., van der Burgt, F. R., & Dijkstra, H. A. (2015). An early warning indicator for atmospheric blocking events using transfer operators. *Chaos: An Interdisciplinary Journal of Nonlinear Science*, 25(3), 036406. <https://doi.org/10.1063/1.4908174>
- Tibaldi, S., & Molteni, F. (1990). On the operational predictability of blocking. *Tellus*, 42(3), 343–365. <https://doi.org/10.1034/j.1600-0870.1990.t01-2-00003.x>
- Vitart, F., Ardilouze, C., Bonet, A., Brookshaw, A., Chen, M., Codorean, C., et al. (2017). The Subseasonal to Seasonal (S2S) prediction project database. *Bulletin of the American Meteorological Society*, 98(1), 163–173. <https://doi.org/10.1175/BAMS-D-16-0017.1>
- Woollings, T., Barriopedro, D., Methven, J., Son, S.-W., Martius, O., Harvey, B., et al. (2018). Blocking and its response to climate change. *Current Climate Change Reports*, 4(3), 287–300. <https://doi.org/10.1007/s40641-018-0108-z>
- Yang, M., Luo, D., Li, C., Yao, Y., Li, X., & Chen, X. (2021). Influence of atmospheric blocking on storm track activity over the North Pacific during boreal winter. *Geophysical Research Letters*, 48(17), e2021GL093863. <https://doi.org/10.1029/2021GL093863>
- Yosinski, J., Clune, J., Bengio, Y., & Lipson, H. (2014). How transferable are features in deep neural networks? *Advances in Neural Information Processing Systems*, 27.
- Zappa, G., Masato, G., Shaffrey, L., Woollings, T., & Hodges, K. (2014a). Linking Northern Hemisphere blocking and storm track biases in the CMIP5 climate models. *Geophysical Research Letters*, 41(1), 135–139. <https://doi.org/10.1002/2013gl058480>
- Zappa, G., Masato, G., Shaffrey, L., Woollings, T., & Hodges, K. (2014b). Linking Northern Hemisphere blocking and storm track biases in the CMIP5 climate models. *Geophysical Research Letters*, 41(1), 135–139. <https://doi.org/10.1002/2013GL058480>
- Zhang, H. (2024). hzhang-math/blockingshaptl: Code\_blocking. Zenodo. <https://doi.org/10.5281/zenodo.13829703>