



# RESEARCH ARTICLE

10.1029/2023MS003809

## Key Points:

- A novelty detector is trained to classify whether an atmospheric profile belongs to a high-resolution model's distribution of profiles
- The detector classifies more profiles as novelties in machine-learning corrected simulations that drift further from a reference simulation
- Using the detector to turn off corrections for novel profiles leads to corrected simulations with more consistent and less biased climates

## Correspondence to:

A. Kwa,  
[annak@allenai.org](mailto:annak@allenai.org)

## Citation:

Sanford, C., Kwa, A., Watt-Meyer, O., Clark, S. K., Brenowitz, N., McGibbon, J., & Bretherton, C. (2023). Improving the reliability of ML-corrected climate models with novelty detection. *Journal of Advances in Modeling Earth Systems*, 15, e2023MS003809. <https://doi.org/10.1029/2023MS003809>

Received 5 MAY 2023  
Accepted 19 OCT 2023

© 2023 The Allen Institute for Artificial Intelligence. Journal of Advances in Modeling Earth Systems 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](#), which permits use, distribution and reproduction in any medium, provided the original work is properly cited.

# Improving the Reliability of ML-Corrected Climate Models With Novelty Detection

Clayton Sanford<sup>1,2</sup> , Anna Kwa<sup>1</sup> , Oliver Watt-Meyer<sup>1</sup> , Spencer K. Clark<sup>1,3</sup> , Noah Brenowitz<sup>4</sup> , Jeremy McGibbon<sup>1</sup> , and Christopher Bretherton<sup>1</sup>

<sup>1</sup>Allen Institute for Artificial Intelligence, Seattle, WA, USA, <sup>2</sup>Department of Computer Science, Columbia University, New York, NY, USA, <sup>3</sup>Geophysical Fluid Dynamics Laboratory, NOAA, Princeton, NJ, USA, <sup>4</sup>NVIDIA Corporation, Santa Clara, CA, USA

**Abstract** Using machine learning (ML) for the online correction of coarse-resolution atmospheric models has proven effective in reducing biases in near-surface temperature and precipitation rate. However, ML corrections often introduce new biases in the upper atmosphere and causes inconsistent model performance across different random seeds. Furthermore, they produce profiles that are outside the distribution of samples used in training, which can interfere with the baseline physics of the atmospheric model and reduce model reliability. This study introduces the use of a novelty detector to mask ML corrections when the atmospheric state is deemed out-of-sample. The novelty detector is trained on profiles of temperature and specific humidity in a semi-supervised fashion using samples from the coarsened reference fine-resolution simulation. The novelty detector responds to particularly biased simulations relative to the reference simulation by categorizing more columns as out-of-sample. Without novelty detection, corrective ML occasionally causes undesirably large climate biases. When coupled to a running year-long coarse-grid simulation, novelty detection deems about 21% of columns to be novelties. This identification reduces the spread in the root-mean-square error (RMSE) of time-mean spatial patterns of surface temperature and precipitation rate across a random seed ensemble. In particular, the random seed with the worst RMSE is improved by up to 60% (depending on the variable) while the best seed maintains its low RMSE. By reducing the variance in quality of ML-corrected climate models, novelty detection offers reliability without compromising prediction quality in atmospheric models.

**Plain Language Summary** Fine-grid global storm-resolving models produce more accurate rainfall and temperature forecasts than coarse-grid climate models, but are too computationally expensive to run for many years. Corrective machine learning (ML) can help coarse-grid climate models act more like fine-grid models, but also makes them more vulnerable to inputs lying outside the range of training data for the ML algorithm. For such “out-of-sample” inputs, the ML may give unreliable results. Using a separate ML scheme, we identify out-of-sample data and disable the ML correction for these cases. We find that this robustly improves the time-mean temperature and precipitation patterns predicted by ML-corrected climate simulations to be 30%–50% better than similar simulations without ML. Incorporating novelty detectors into ML-corrected simulations can improve their prediction skill by helping them avoid drifting into “out-of-sample” states.

## 1. Introduction

Accurate, reliable climate models are essential for projecting climate change and its impacts. To explore a range of scenarios and account for natural climate variability, climate models must also be computationally efficient. This is typically achieved in the atmospheric model by using relatively coarse grid resolutions (typically between 50 and 200 km) and representing processes that operate at finer spatial scales by somewhat empirical human-designed “subgrid parameterizations.”

The use of machine learning in atmospheric modeling has taken various forms, including emulating existing physical parameterizations (e.g., Chantry et al., 2021; Krasnopolsky et al., 2010), replacing physics parameterizations by learning from a high-resolution model (e.g., Brenowitz & Bretherton, 2019; Rasp et al., 2018; Wang et al., 2022; Yuval & O’Gorman, 2020), or using ML for online correction of a complete atmospheric model (Bretherton et al., 2022; Chen et al., 2022; Clark et al., 2022; Kwa et al., 2022; Watt-Meyer et al., 2021). Here we will focus on the latter strategy.

Previous works (Brenowitz & Bretherton, 2019; Rasp et al., 2018; Watt-Meyer et al., 2021; Yuval & O’Gorman, 2020) have suggested that correcting or augmenting physics-based climate models with machine learning (ML) can improve weather forecast skill and reduce climate biases. However, ML-augmented models can be susceptible to instabilities (Brenowitz, Beucler, et al., 2020), and their performance when coupled to the atmospheric model can be sensitive to subtle ML training differences, such as random seed selection (Brenowitz, Henn, et al., 2020; Wang et al., 2022).

This study draws on the idea of using a compound parameterization (Krasnopolsky et al., 2008; Song et al., 2021) to mask ML models with high uncertainty. Specifically, we train a novelty detection algorithm (Hodge & Austin, 2004) and use it at each timestep of a coarse-grid simulation to mask ML corrections when the column atmospheric state is determined to be outside the distribution of the data set used to train the ML model estimating online corrections. Our approach adds robustness to past approaches (specifically Kwa et al., 2022) while consistently improving temperature and precipitation bias metrics. A preliminary version of this study was presented in C. H. Sanford et al. (2022); an important but unrelated software bug fix and some changes in configuration led to substantial changes in interpretation of the effects of the novelty detector, as discussed in Section 2.5.

We model the atmosphere as a discretized system of partial differential equations. The atmospheric state is modeled as  $X = (x_1, \dots, x_N) \in \mathbb{R}^{N \times d}$ , a three-dimensional grid of  $N$  latitude/longitude coordinates with  $d$ -dimensional column vectors concatenating the vertical profiles of gridpoint values of air temperature, specific humidity, winds and other fields. In a “baseline” model with no added ML corrections, the state of a particular column  $x_i \in \mathbb{R}^d$  evolves over time as

$$\frac{dx_i}{dt} = f_i(X, t) \quad (1)$$

for some fixed  $f_i$  derived from physically based assumptions.

The number of grid columns  $N$  scales with the inverse square of the desired grid spacing; large  $N$  (a fine grid) typically yields more accurate average estimates of the temperature, humidity, and precipitable water in the atmosphere, at the cost of computational efficiency. While accuracy penalties due to poor grid resolution are expected for small  $N$ , coarse-grid simulations are also biased by imperfect representations of subgrid-scale processes like thunderstorms and cloud radiative effects (Woelfle et al., 2018; Zhang & Wang, 2006). ML is an appealing way to de-bias this coarse climate model by predicting and compensating for its error. The ML-corrected model can be written

$$\frac{dx_i}{dt} = f_i(X, t) + g(x_i, \varphi_i; \theta), \quad (2)$$

where  $g(\cdot; \theta) : \mathbb{R}^{d+3} \rightarrow \mathbb{R}^d$  is a learned function with parameters  $\theta$  that predicts corrective tendencies from the column,  $x_i \in \mathbb{R}^d$ , and its insolation, surface elevation, and latitude  $\varphi_i \in \mathbb{R}^3$ . The ML correction enables the baseline to better approximate a reference fine-grid model while maintaining the underlying physics as the core of the modeling approach (Brenowitz & Bretherton, 2019; Watt-Meyer et al., 2021).

The ML model is trained and evaluated “offline” by generating predictions over single timesteps from their corresponding input state columns in a fixed data set. In “online” application, the model is coupled to the components of the coarse-grid atmosphere model as the year-long forecast is simulated. While ML-based models frequently improve overall error, these models—especially deep neural networks—are often not robust, meaning they perform poorly for out-of-sample data. In online application, where predictions are fed back into the model, the corrective ML can induce errors in the overall simulation that accumulate in time, creating large systematic biases and instabilities (Brenowitz, Henn, et al., 2020).

This motivated us to employ semi-supervised novelty detection to predict when a column  $x_i$  belongs to the training distribution of  $g$  and suppress the tendencies of the ML model if not. This paper shows that strategy can substantially improve the model stability and climate accuracy. With novelty detection, our model has the form

$$\frac{dx_i}{dt} = f_i(X, t) + \eta(x_i; \rho)g(x_i, \varphi_i; \theta), \quad (3)$$

for a novelty detector  $\eta(\cdot; \rho) : \mathbb{R}^d \rightarrow [0, 1]$  with parameter vector  $\rho$ .

## 2. Methodology

### 2.1. Data Set

We train the ML tendency correction  $g(\cdot; \theta)$  offline as described by Kwa et al. (2022). The training samples  $((x_1, \varphi_1), y_1), \dots, ((x_n, \varphi_n), y_n)$  consist of input features and target nudging tendencies  $y_i$  (described below) for a set of atmospheric columns sampled at time steps of a nudged coarse model simulation. A neural net with parameters  $\theta$  is trained to make  $g$  best match the nudging tendencies.

The nudged coarse model simulation is constructed to track the evolution of a reference fine-grid no-ML climate model simulation, averaged to the coarse grid cells. Symbolically, the atmospheric state in this reference simulation is denoted  $X_{\text{fine}}^{(1)}, \dots, X_{\text{fine}}^{(T)} \in \mathbb{R}^{N \times d}$ .

To nudge the coarse simulation to this fine-grid reference, we add a relaxation term- referred to as a nudging tendency- to the coarse-grid model of the form

$$y_i := \frac{X_{\text{fine},i} - X_i}{\tau}, \quad (4)$$

with a specified nudging timescale  $\tau = 3$  hr. By construction, the time-evolving atmospheric state  $X^{(1)}, \dots, X^{(T)}$  of this nudged run is approximately (but not exactly) the same as in the fine-grid reference. The  $y_i$  are the nudging tendencies that we learn; we denote the  $N \times d$  arrays of their values at each time as  $Y^{(1)}, \dots, Y^{(T)}$ .

The coarse-grid model  $f_i$  is the same as used in Bretherton et al. (2022) and Kwa et al. (2022). We use a version of NOAA's FV3GFS global weather forecast model (Zhou et al., 2019) with a C48 cubed-sphere grid of approximately 200 km horizontal grid spacing (Putman & Lin, 2007). In this grid, the Earth is divided into 6 square tiles with a 48-by-48 grid imposed on each, for  $N = 6 \cdot 48^2$  grid columns. This model has 79 vertical levels between the surface and the top of the atmosphere. Time-varying sea surface temperature and sea ice are prescribed using a coarsened version of the same fields from the fine-grid reference.

Our fine-grid reference simulation is the same as used in Kwa et al. (2022). It is created using the X-SHIELD model, a modified configuration of FV3GFS on a C3072 (approximately 3 km) cubed-sphere grid with 79 vertical model levels (Cheng et al., 2022). The FV3GFS convective gravity wave drag and deep cumulus parametrization schemes are disabled in the fine-grid model, while the shallow cumulus convection scheme is active. We used a year of three-hourly reference model output coarsened to the C48 grid by horizontal pressure-level averaging (Bretherton et al., 2022).

Samples are collected from a year-long nudged coarse-grid simulation; the state and nudging tendencies are saved every 3 hr. After dividing this data into interleaved time blocks for the train/test split and randomly subsampling down to 15% of the columns in each timestep, we are left with  $n = 2834611$  training samples spanning 2020-01-19 through 2021-01-17.

The same data set  $\mathcal{D}_x = \{x_i \in \mathbb{R}^d : i \in [n]\}$  is used to train the novelty detector  $\eta(\cdot; \rho)$ . The nudging tendencies  $y_i$  are omitted, as the novelty detection procedure requires no labels.

### 2.2. ML-Corrected Climate Models and Data

The novelty detection procedure does not affect the training of the neural nets used to predict the nudging tendencies. We consider two such corrective ML models:  $g_{Tq}$  and  $g_{Tquv}$ :

- $g_{Tq}$  corrects vertical columns of air temperature  $T$  and specific humidity  $q$  tendencies, but does not correct winds. That is,  $x_i$  is a  $d = (2 \cdot 79)$ -dimensional vector with 79 temperature and 79 humidity coordinates, each corresponding to an atmospheric model level.
- $g_{Tquv}$  also corrects tendencies of the horizontal wind components  $(u, v)$  at each level, making  $x_i$  a  $d = (4 \cdot 79)$ -dimensional vector.

$g_{Tq}(\cdot; \theta) : \mathbb{R}^{158} \times \mathbb{R}^3 \rightarrow \mathbb{R}^{158}$  predicts the vector of temperature and humidity nudging tendencies  $y_i$  from the temperature and humidity profiles  $x_i$ , as well as the insolation, surface elevation, and latitude of the corresponding cell  $(\varphi_i)$ . Hyperparameters for the corrective ML models were selected after performing a sweep optimized on single-timestep validation loss. We represent  $g_{Tq}(\cdot; \theta)$  as a three-layer dense multi-layer perceptron of width 419.

The loss is measured by the mean absolute error (MAE) with  $L_2$  kernel regularization of strength  $10^{-4}$ . We found that models trained with MAE loss were less prone to instabilities and drifts in online simulations than those trained with mean squared error loss. We train the model with the Adam optimizer for 500 epochs using a fixed learning rate of 0.00014 and a batch size of 512 samples.

On the other hand,  $g_{\text{Tquv}}(\cdot; \theta) : \mathbb{R}^{316} \times \mathbb{R}^3 \rightarrow \mathbb{R}^{316}$  is defined as the concatenation of two learned functions for input  $x = (x_{\text{Tq}}, x_{\text{uv}}) \in \mathbb{R}^{158} \times \mathbb{R}^{158}$ :

$$g_{\text{Tquv}}(x, \varphi; \theta) = (g_{\text{Tq}}(x_{\text{Tq}}, \varphi; \theta_{\text{Tq}}), g_{\text{uv}}(x_{\text{Tq}}, x_{\text{uv}}, \varphi; \theta_{\text{uv}})) \quad (5)$$

where  $g_{\text{Tq}}(\cdot; \theta_{\text{Tq}})$  is the same as the aforementioned model.  $g_{\text{uv}}(\cdot; \theta_{\text{uv}}) : \mathbb{R}^{316} \times \mathbb{R}^3 \rightarrow \mathbb{R}^{158}$  is separately trained to infer wind nudging tendencies from temperatures, humidities, and horizontal winds. Besides the different input dimension,  $g_{\text{uv}}(\cdot; \theta_{\text{uv}})$  is otherwise structured and trained identically to the other model.

### 2.2.1. Fixed Vertically Flipped Application of Corrective Wind Tendencies

Kwa et al. (2022) obtained better prognostic simulations using  $g_{\text{Tq}}$  than by also adding wind tendency correction  $g_{\text{Tquv}}$ . We have since found this was due to our inadvertently applying the learned wind tendency correction in each column upside-down, such that the correction of the lowest level 79 was applied on the highest level 1, and vice versa. This configuration error, which also affected the wind-corrected simulations discussed by Watt-Meyer et al. (2021), Bretherton et al. (2022), and Clark et al. (2022), arose because FV3GFS uses opposite vertical indexing of grid levels in the physical parameterizations and dynamical core. After fixing this error, including corrective wind tendencies no longer leads to numerical instability, and most metrics of 3–7 days weather skill (e.g., RMSEs of 850 hPa temperature) are significantly improved.

We tested both the erroneous and the fixed corrective approaches to trace through the effects of rectifying the error caused by vertical flipping. Both corrective approaches benefit from the inclusion of novelty detection.

### 2.3. Novelty Detection

Novelty detection is a well-studied semi-supervised learning problem about estimating the support of a data set using only positive examples (Hodge & Austin, 2004). Most novelty detection algorithms predict whether a new sample  $x$  (which does not appear in the finite-size training data set) belongs in the support of the training data set, which is a subset  $S$  of the input space such that all columns  $x \in S$  have some positive probability of being drawn by the probability distribution used to generate the training set. Lacking access to any explicit characterization of this training distribution, novelty detectors estimate the set  $S$  statistically, using the training features without labels (i.e., positive examples). We frame the problem as novelty detection rather than outlier detection (an unsupervised problem with mixture of in-distribution and out-of-distribution samples) or standard two-class supervised classification, because we have no data set of representative out-of-distribution samples and constructing such a data set would introduce additional model-dependence into this process.

In our work, the novelty detector  $\eta$  predicts an estimate of the support of the training data  $\hat{S}$ , which we use to mask the ML-predicted corrections of input  $x$  when  $x \notin \hat{S}$ . Specifically, if a column is determined to not be a novelty (i.e.,  $x \in \hat{S}$ ), then we let  $\eta(x; \rho) = 1$  (recall Equation 3) to take full advantage of the learned correction  $g(x, \varphi; \theta)$ ; otherwise, we ignore  $g(x, \varphi; \theta)$  by setting  $\eta(x; \rho) = 0$ .

There are many known approaches to novelty detection, including local-outlier factor (Breunig et al., 2000),  $k$ -means clustering (Nairac et al., 1999), and minimum-volume ellipsoid estimation (Van Aelst & Rousseeuw, 2009). Our exploratory work considers two approaches: a simple “min-max” novelty detector and a one-class support vector machine (OCSVM). For each of these we consider novelty detectors  $\eta_{\text{T}}$  with 79-dimensional temperature vectors as input and  $\eta_{\text{Tq}}$  with 158-dimensional combined temperature and specific humidity vectors.

We did not consider novelty detectors with wind inputs. Adding more inputs to the OCSVM classifier requires further hyperparameter tuning (see Section 2.3.3) to keep the evaluation time low enough to be useable within prognostic simulations; we therefore limit the scope of this work to out-of-sample detection on temperature and specific humidity fields.

### 2.3.1. Naive “Min-Max” Novelty Detector

The min-max novelty detector considers the smallest axis-aligned hyper-rectangle that contains all training samples and categorizes any sample outside the rectangle as a novelty:

$$\eta_{\text{min-max}}(x; (x_{\min}, x_{\max})) = \begin{cases} 1 & \text{if } x_k \in [x_{\min,k}, x_{\max,k}] \forall k \in [d], \\ 0 & \text{otherwise,} \end{cases} \quad (6)$$

for  $x_{\min,k} = \min_{i \in \mathcal{I}} x_{i,k}^{(n)}$  and  $x_{\max,k} = \max_{i \in \mathcal{I}} x_{i,k}^{(n)}$  as the minimum and maximum over the training data of the  $k$ th feature. While efficient, this novelty detector cannot identify irregular correlations between input features that nevertheless lie within the bounding box.

### 2.3.2. One-Class Support Vector Machine (OCSVM)

The one-class SVM algorithm of Schölkopf et al. (2001) repurposes the SVM classification algorithm to estimate the support of a distribution by finding the maximum-margin hyperplane separating training samples from the origin. The OCSVM has been applied to novelty detection for genomics (Sommer et al., 2017), video footage (Amraee et al., 2018), propulsion systems (Tan et al., 2019), and the internet of things (Yang et al., 2021).

We normalize each input  $x_i$  and utilize the kernel trick, lifting it to the infinite-dimensional feature space  $\phi(x_i)$  corresponding to the radial basis function (RBF) kernel  $\kappa_\gamma(x, x') = \exp(-\gamma \|x - x'\|_2^2)$ . We parameterize the novelty detector with  $\rho = (\alpha, \xi, \gamma)$  in its dual form,

$$\eta_{\text{OCSVM}}(x; (\alpha, \xi, \gamma)) = \begin{cases} 1 & \text{if } \sum_{i=1}^n \alpha_i \kappa_\gamma(x, x_i) \geq \xi, \\ 0 & \text{otherwise.} \end{cases} \quad (7)$$

The sensitivity of the novelty detector can be adjusted by choosing a cutoff  $\xi > 0$ . The learnable real-valued weights  $\alpha_i \geq 0$  quantify the influence of each training sample  $x_i$  on predictions on any new sample  $x$ . By directly weighting the kernel  $\kappa_\gamma(x, x_i)$ , a large  $\alpha_i$  indicates that the proximity of  $x$  to  $x_i$  is highly salient to  $\eta_{\text{OCSVM}}$ 's prediction of whether  $x_i$  belongs in the support. If  $\alpha_i = 0$ , then  $\kappa_\gamma(x, x_i)$  is irrelevant to the prediction and need not be computed. The goal is to find a relatively small subsample of training samples  $x_i$  with nonzero weights  $\alpha_i$ , known as the support vectors, that can be used to confidently and efficiently assess whether  $x$  is out of sample.

The weights  $\alpha_i$  are learned by solving a convex optimization problem based on the training data. The number of nonzero weights  $\alpha_i$  depends on the sensitivity  $\gamma$  and a regularization parameter  $\nu$ . To obtain a robust and computationally efficient novelty detector, for a given  $\gamma$  we choose  $\nu$  to ensure the number of support vectors is on the order of at most  $10^4$ , less than 0.5% of the training data sample. Smaller values of  $\gamma$  correspond to novelty detectors with highly smoothed support estimations that may be larger than necessary, while large  $\gamma$  provides a smaller and perhaps more topologically complex region.

### 2.3.3. Parameterization of the One-Class Support Vector Machine

For the main results presented in 3, we chose  $\gamma = 4/79$ ,  $\nu = 10^{-4}$ , and  $\xi = 0.12$  as our OCSVM model parameters. Here, we discuss the process of OCSVM parameter selection and the resulting trade-offs. Specifically,  $\gamma$  was set to moderate a bias-variance trade-off (although a wide variety of choices produce similar results),  $\nu$  guarantees that the costs of computing the outcome of the novelty detector are dominated by the other steps of the simulation, and  $\xi$  was set to tune the classification boundary.

The RBF kernel in our OCSVM can be tuned to trade off bias and variance with its inverse-radius parameter  $\gamma$ . A large choice of  $\gamma$  ensures that  $\kappa_\gamma(x, x')$  only has non-negligible output if  $x$  is extremely close to  $x'$ , while smaller  $\gamma$  selections cause a large “ball” of  $x$  around  $x'$  to all have  $\kappa_\gamma(x, x') \approx 1$ . Choosing large  $\gamma$  makes for a more expressive classifier that can be used to fit any training data, but raises the risk of classifying many “holes” in between training data samples as out-of-sample. A smaller  $\gamma$  imposes a smoothing effect on the learned classifier. The default SVM setting in scikit-learn is  $\gamma = \frac{1}{\text{\# features}} = \frac{1}{2.79}$ . For our application, we find that a larger choice of  $\gamma$  tends to produce better outcomes and focus our study on four choices:  $\gamma \in \left\{ \frac{1}{79}, \frac{2}{79}, \frac{4}{79}, \frac{8}{79} \right\}$ .

**Table 1**  
One-Class SVM Parameterizations

| $\gamma$       | $\nu$                               | # SVs         | $\xi_{0.25}$        | $\xi_{0.5}$         | $\xi_{0.75}$        | $\xi_{0.95}$        | $\xi_{0.99}$        |
|----------------|-------------------------------------|---------------|---------------------|---------------------|---------------------|---------------------|---------------------|
| $\frac{1}{79}$ | $5 \cdot 10^{-3}$                   | 14,365        | 351                 | 321                 | 289                 | 227                 | 153                 |
| $\frac{2}{79}$ | $5 \cdot 10^{-3}$                   | 15,029        | 80                  | 70                  | 60                  | 42                  | 22                  |
| $\frac{4}{79}$ | <b><math>1 \cdot 10^{-4}</math></b> | <b>16,030</b> | 0.18                | 0.15                | <b>0.12</b>         | 0.065               | 0.023               |
| $\frac{8}{79}$ | $4 \cdot 10^{-6}$                   | 12,152        | $5.9 \cdot 10^{-4}$ | $4.4 \cdot 10^{-4}$ | $2.8 \cdot 10^{-4}$ | $9.3 \cdot 10^{-5}$ | $1.7 \cdot 10^{-5}$ |

*Note.* For each kernel radius  $\gamma$ , we select a regularization parameter  $\nu$  in order to constrain the number of support vectors to roughly 10,000 for computational efficiency, which is in turn used to train a parameter vector  $\alpha$ . Five cutoffs  $\xi$  are identified to adjust the conservatism of the model:  $\xi_p$  is chosen to ensure that a  $p$  fraction of the training data set is categorized as in-distribution, that is,  $\eta_{\text{OCSVM}}(x; (\alpha, \xi_p, \gamma)) = 1$ . Bold values indicate the OCSVM parameters used in the main results.

The scikit-learn implementation of an OCSVM uses a regularization parameter  $\nu$  in the training procedure to trade off classification accuracy and model simplicity when learning weights  $\alpha \in [0, 1]^n$  (Schölkopf et al., 2000).  $\nu$  does so by regulating the number of allowable support vectors, which are samples  $x_i$  that have respective weight  $\alpha_i > 0$ . A looser bound on support vectors in turn scales the computational cost of each application of the OCSVM. Choosing a large value of  $\nu$  puts a greater premium on categorizing every sample correctly by using more support vectors. Here, we use a parameter search to choose a  $\nu$  for each  $\gamma$  that results in roughly  $10^4$  support vectors.

Finally, the cutoff  $\xi$  affects the sensitivity of the learned novelty detector. A large choice of  $\xi$  causes an aggressive detector that categorizes a large number of samples as novelties (and hence, frequently disables the ML-corrected tendencies), while a small  $\xi$  classifies more samples as in-distribution. We use the maximum score observed in the training data,  $\xi = 0.12$ , which classifies none of the training data and an acceptably small 2.6% of a withheld test set of reference data as out-of-sample.

Section 4 investigates the dependence of simulations' accuracy metrics on several choices of the sensitivity  $\gamma$  and cutoff  $\xi$ . We calibrate the sensitivity by drawing samples from a full year of an ML-corrected run and choosing a cutoff  $\xi_p$  such that a fraction  $p$  of the given data are categorized as in-distribution; a larger choice of  $p$  results in a smaller  $\xi_p$ . For the sensitivity study in Section 4, we consider the corresponding  $\xi_p$  choices for each  $\gamma$  for  $p \in \{0.25, 0.5, 0.75, 0.95, 0.99\}$ . The value of  $\xi = 0.12$  used in our results corresponds to  $p = 0.75$  when evaluated on the ML-corrected run. In Table 1, we give the respective choices of  $\nu$  and  $\xi_p$  for each  $\gamma$ .

## 2.4. Computing Scalar Metrics

We measure the success of a coarse-grid simulated run by computing the root mean-square error (RMSE) of time-averaged quantities (850 hPa and 200 hPa temperature, surface precipitation, total precipitable water) with respect to those same quantities for the coarsened fine-grid run. We compute the RMSE of the time-averaged field  $s$  as follows:

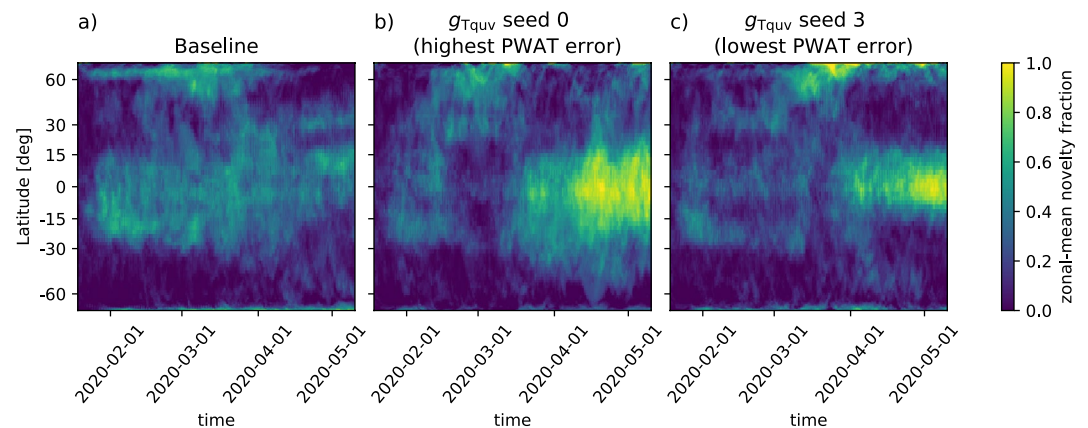
$$\text{RMSE}(s) = \sqrt{\sum_{i=1}^N a_i \left( \frac{1}{T} \sum_{t=1}^T \left( \hat{s}_i^{(t)} - s_{\text{fine},i}^{(t)} \right)^2 \right)}, \quad (8)$$

where  $\hat{s}_i^{(t)}$  and  $s_{\text{fine},i}^{(t)}$  denote the field value at grid cell  $i \in [N]$  and time  $t \in [T]$  in our coarse-grid and the reference fine-grid simulations respectively, and  $a_i$  are the normalized area weights of grid cells.

## 2.5. Methodological Updates Versus Sanford et al. (2022)

We made two important methodological updates in this study compared to a similar recent work on which it is based (C. H. Sanford et al., 2022). First, we fixed the previous error (see Section 2.2.1), discovered after that earlier work, where the ML wind tendencies in each grid column were applied with inverted vertical indexing during online simulations. The second change is related to the application of the ML corrections  $g_{Tq}$  and  $g_{Tquv}$





**Figure 1.** Zonal-mean fraction of novelties detected by the  $\eta_{Tq,OC SVM}$  novelty detector over the first 16 weeks of the (a) baseline, (b) seed 0  $g_{Tquv}$  and (c) seed 3  $g_{Tquv}$  simulations.

in the upper atmosphere. C. H. Sanford et al. (2022) followed the approach of Kwa et al. (2022), in which the ML-predicted tendencies in the top three model layers were not applied as corrections. The rationale was that the sponge layer differences between low and high resolution models was a process we did not wish to correct, and there were relatively large magnitude nudging tendencies at these levels. In this study, we use a more aggressive tapering in which the ML-predicted outputs are tapered to zero throughout the uppermost 25 model levels using an exponential decay, as in Equation 6 of Clark et al. (2022). This improves the simulation of lower atmospheric air temperatures, and more importantly, helps prevent large upper atmospheric temperature drifts when using ML corrections of horizontal winds. Both of these changes improve the ML-corrected simulations described by Equation 2 and impose a higher bar for the novelty detection to add value.

### 3. Results

#### 3.1. Offline Application of Novelty Detection

Before integrating a novelty detector into online simulations with an ML-corrected climate model, we test it offline on data produced by the preexisting simulations. We compare the frequency of offline novelty detection for data sets generated from the first 16 weeks of three C48 simulations—a no-ML baseline model simulation and two  $g_{Tquv}$ -corrected simulations that differ only in the random initial seed used in training the  $g_{Tquv}$  models. The  $g_{Tquv}$  seed 0 run has the largest yearly mean precipitable water RMSE ( $4.4 \text{ kg/m}^2$ ) across a set of four  $g_{Tquv}$  simulations, while the seed 3 run has the smallest ( $2.4 \text{ kg/m}^2$ ), slightly smaller than that of the baseline run ( $2.7 \text{ kg/m}^2$ ). Feedback loops between less reliable ML corrections and out-of-sample column states may exacerbate mean-state drifts, showing up as locally higher offline novelty fractions. We selected the two ML-corrected runs with the highest and lowest precipitable water biases in order to test this hypothesis. The baseline simulation tests the extent to which mean-state biases developing in a conventional climate model lead to detectable novelties.

Figure 1 focuses on the first 16 simulated weeks of the simulation to make the drifts into out-of-sample states more visible. Within a few days, the baseline model moistens relative to the reference model until it generates enough clouds and precipitation to balance surface evaporation, after which it settles into a new, slightly biased equilibrium in which about 25% of the columns are flagged as novelties.

Initially, the seed 0 and seed 3  $g_{Tquv}$  corrections both have the intended effect of keeping the global state closer to the fine-grid reference distribution. These ML-corrected  $g_{Tquv}$  runs have lower global novelty fractions than the baseline over the first 2 months, particularly in the tropics. However, from March onward, the novelty fraction in the baseline tropics plateaus, while both  $g_{Tquv}$  simulations continue to drift farther out-of-sample in the tropics.

By the end of the 16 weeks shown in Figure 1, the “highest PWAT error” seed 0  $g_{Tquv}$  simulation has roughly twice as many out-of-sample columns compared to the baseline and “lowest PWAT error” seed 3 runs. This demonstrates that suboptimal ML corrections (as in the seed 0  $g_{Tquv}$  model) can indeed push the state further out

**Table 2**

*The Root Mean-Square Error Scores of Time-Averaged Metrics and Novelty Detection Rates for Year-Long Simulations*

| Run                               | % Novelty  | T200 (K)           | T850 (K)           | SP (mm/day)        | PWAT (kg/m <sup>2</sup> ) |
|-----------------------------------|------------|--------------------|--------------------|--------------------|---------------------------|
| Baseline                          | -          | 2.48               | 2.09               | 1.78               | 2.79                      |
| $g_{Tq}$                          | -          | 2.50 (0.40)        | 1.97 (0.08)        | 1.52 (0.07)        | 3.97 (0.29)               |
| $g_{Tquv}$                        | -          | 3.30 (0.49)        | 1.31 (0.14)        | 1.40 (0.12)        | 3.40 (0.73)               |
| $g_{Tquv}, \eta_{T, \min - \max}$ | 0.6 (0.3)  | 3.04 (0.65)        | <b>1.29 (0.06)</b> | 1.36 (0.07)        | 3.28 (0.72)               |
| $g_{Tquv}, \eta_{T, OCSVM}$       | 5.0 (1.0)  | 2.84 (0.49)        | 1.38 (0.09)        | 1.37 (0.08)        | 3.36 (0.83)               |
| $g_{Tquv}, \eta_{Tq, OCSVM}$      | 20.6 (4.8) | <b>1.24 (0.05)</b> | 1.30 (0.08)        | <b>1.29 (0.07)</b> | <b>2.38 (0.37)</b>        |

*Note.* Values for ML-corrected runs are the mean, with standard deviation in parentheses, across the four random seeds. The “% Novelty” column represents the percent of columns over the simulated year which were classified as out-of-sample and did not receive ML corrections. Metrics are 200- and 850-hPa temperature (T200, T850), surface precipitation rate (SP) and precipitable water (PWAT). For each metric, the run with the lowest RMSE is bolded.

of the training set distribution, setting the stage for less reliable ML corrections that further exacerbate climate drifts. The higher rate of novelty detection in the tropics and extratropics in the seed 0 simulation is correlated with higher moist biases (not shown) in those regions than in the seed 3 run. This should not be interpreted as the sole physical driver behind the out-of-sample drifts though, as the  $\eta_{Tq, OCSVM}$  novelty detector uses the full column profiles of air temperature and specific humidity in its classification.

### 3.2. Online Novelty Detection Improves Temperature and Precipitation Predictions

We assess the utility of the novelty detectors by incorporating  $\eta(\cdot; \rho)$  into the coarse grid model and numerically simulating Equation 3 for 1 year. We compare the predicted atmospheric states  $\hat{x}_i$  to  $x_{fine, i}$  using the RMSE of four time-averaged diagnostics calculated using Equation 8: air temperatures at pressures of 200 hPa and 850 hPa (T200, T850) representative of the lower and upper troposphere, surface precipitation rate (SP) (Current climate models make less consistent predictions of regional shifts in precipitation than of surface temperatures; contrast Sections B.2.1 and B.3.1 of IPCC (2021).), and precipitable water (PWAT) (PWAT is the total mass of water contained in a vertical atmospheric column per cross-sectional area and is highly correlated with the regional precipitation rate (Bretherton et al., 2004).).

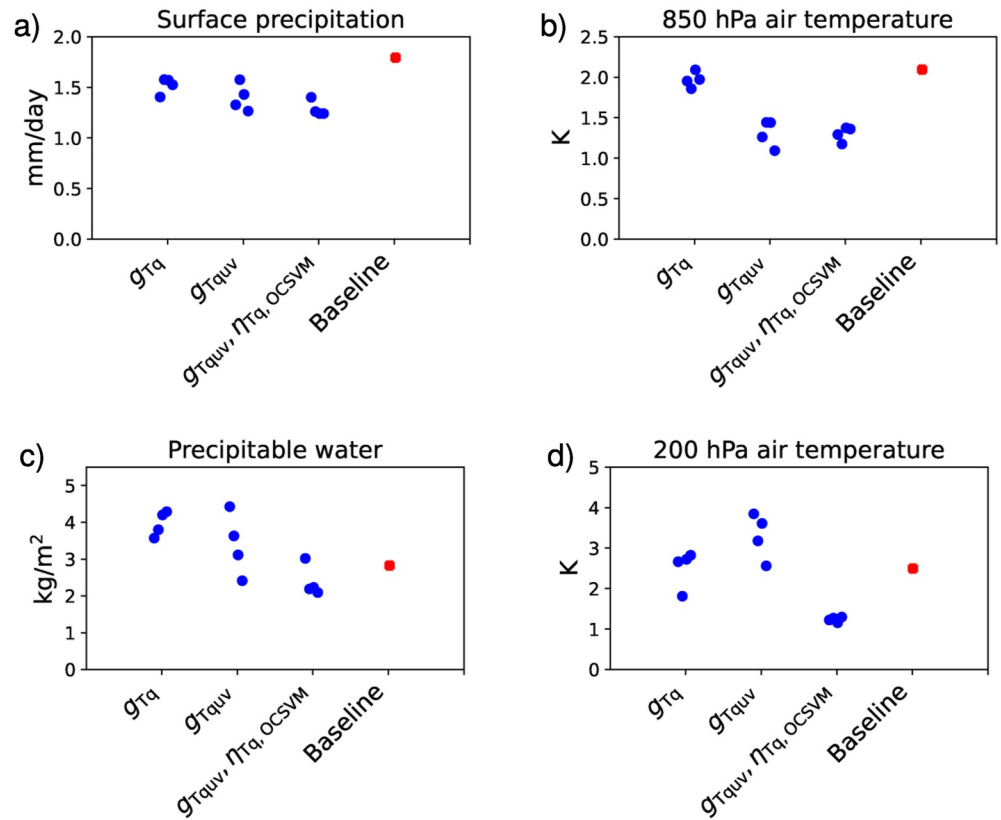
Table 2 compares the performance of six global simulations. The first is the no-ML baseline simulation; the next two are ML-corrected runs without and with wind tendency corrections; and the remaining three simulations use  $g_{Tquv}$  corrections and include novelty detection from Equation 3—these differ in the choice of novelty detector  $\eta$  and its inputs. The  $\eta_{Tq}$  OCSVM uses the same parameter choices as for the offline comparisons. For the  $\eta_T$  OCSVM, which uses fewer inputs, we use the same  $\gamma = 4/79$  and  $\nu = 10^{-4}$  but readjust the cutoff  $\xi$  to 2.02 to the minimum needed to suppress  $T$ -only novelties within the training data set. For all the configurations except the baseline, we perform an ensemble of simulations using four identically trained ML-correction models  $g$  initialized with different random seeds. These are identical to the ML-corrective models used in Kwa et al. (2022) in order to enable direct comparison to the year-long simulations in that previous work.

Without a novelty detector, the conclusions for the  $g_{Tq}$  model (ML-corrected temperature and humidity tendencies only) are similar to Kwa et al. (2022). The metrics (second row in Table 2) are 10%–20% better than for the baseline model, except for the PWAT RMSE which worsens. Adding corrective ML for winds (third row in Table 2) significantly improves the 850 hPa air temperature errors (ensemble-mean RMSE decreases from 1.97 to 1.31 K), somewhat improves SP and PWAT, but substantially worsens the T200 RMSE.

The min-max novelty detector (fourth row in Table 2) slightly improves the RMSEs but has limited impact since it activates only rarely (in 0.6% of atmospheric columns, as shown in the second column of the table). This indicates the importance of bounding the data distribution more tightly than a high-dimensional box. The  $\eta_{T, OCSVM}$  novelty detector classifies a higher fraction of columns as novelties (5%) than the min-max detector, but the overall RMSE for the  $g_{Tquv}, \eta_{T, OCSVM}$  simulations are mostly on par with the  $g_{Tquv}$  results without novelty detection, with the exception of further improvements in T200 RMSE.

The  $\eta_{Tq, OCSVM}$  novelty detector, on the other hand, improves 200 hPa air temperature, surface precipitation, and precipitable water RMSEs by 62%, 8%, and 30% respectively, compared to the  $g_{Tquv}$  simulations without novelty





**Figure 2.** Root mean-square error of time-mean fields in groups of ML-corrected simulations and the baseline prognostic run. Each group of four blue points shows a range of results across four randomly seeded corrective-ML models. The same randomly seeded  $g_{Tq}$  models are used in all ML-corrected groups. The same four  $g_{quv}$  models are used in both the  $g_{Tquv}$  and  $g_{Tquv}, \eta_{Tq}, OCSVM$  groups.

detection. To achieve these improvements, the OCSVM novelty detectors activate in 21% of all atmospheric columns, averaged over the course of the year-long simulations. If compared to the same 16 weeks time period as the offline analysis of  $g_{Tquv}$  runs without novelty detection in Section 3.1, online novelty detection reduces the novelty fraction in ML-corrected runs by roughly half. In summary, suppressing ML corrections to columns with atypical temperature and specific humidity profiles helps keep the  $g_{Tquv}$ -corrected model within the envelope of its training data, where it is skillful in reducing temperature and humidity biases.

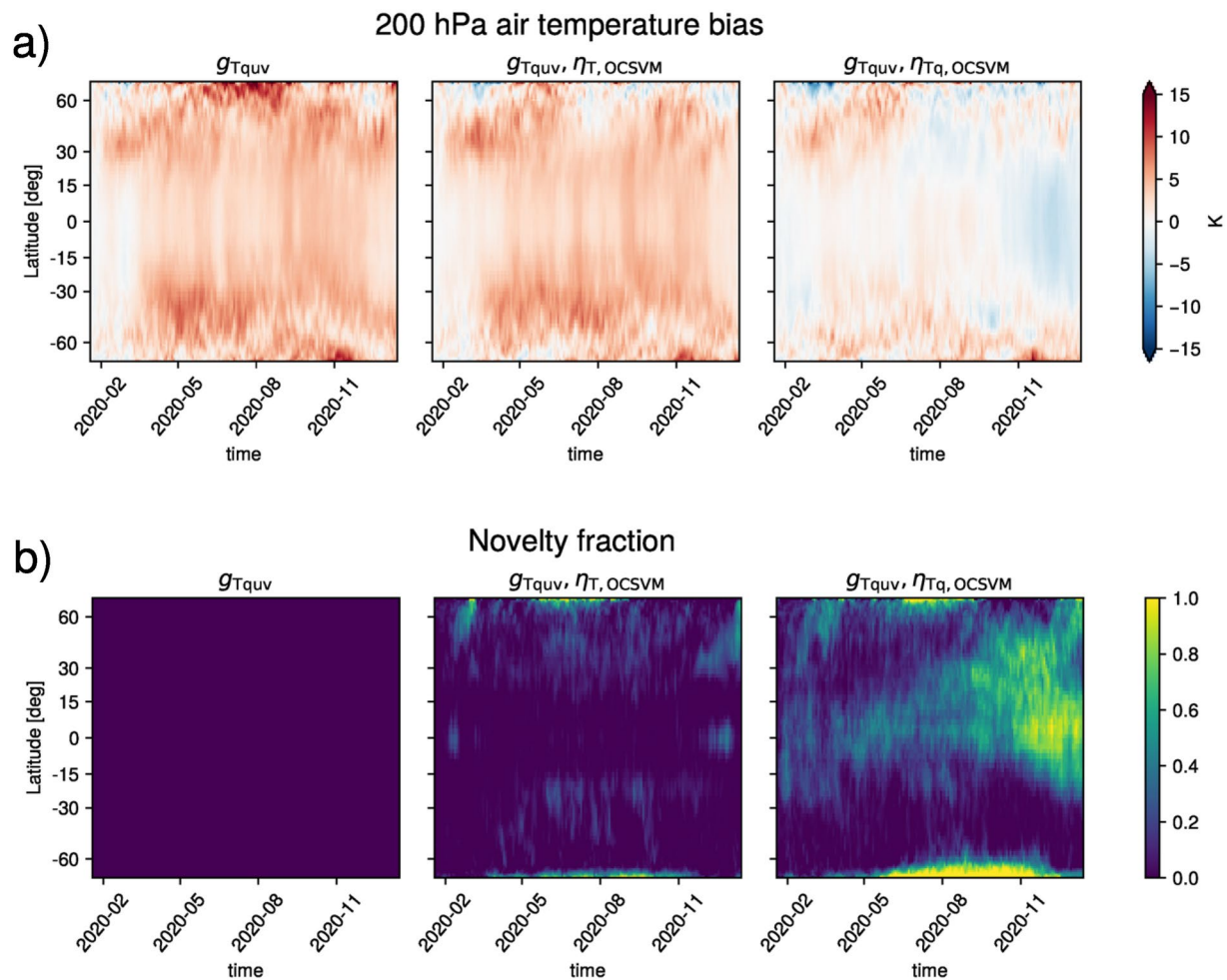
Figure 2 shows the RMSE of time-mean surface precipitation, 200 hPa and 850 hPa temperature, and precipitable water across individual ensemble members of simulations using  $g_{Tq}$ ,  $g_{Tquv}$ , and  $g_{Tquv}, \eta_{Tq}, OCSVM$ . This illustrates that the  $\eta_{Tq}, OCSVM$  novelty detection substantially reduces the variance in skill across the ML-corrected runs (also demonstrated by the standard deviations reported in parentheses in Table 2), especially for precipitable water and 200 hPa temperature. The novelty detection reduces variance and improves the overall ensemble skill by bringing the worst-performing  $g_{Tquv}$  seeds closer in line with the better performers.

### 3.3. Improvements for a Particular ML-Corrected Simulation

In this subsection, the ML-corrected simulation results are shown just for the worst  $g_{Tquv}$  seed (0), to provide a clear illustration of how novelty detection especially benefits poorly performing prognostic runs. This seed's  $g_{Tquv}, \eta_{Tq}, OCSVM$  simulation had a novelty fraction of 24.3%, slightly higher than the ensemble mean of 20.6%.

#### 3.3.1. Zonal-Mean Biases

Figure 3 compares the time evolution of zonal-mean 200 hPa air temperature biases in three ML-corrected year-long simulations:  $g_{Tquv}$  without novelty detection, and two simulations with novelty detectors  $\eta_{Tq}, OCSVM$  and  $\eta_{Tq}, OCSVM$  that use different feature sets.

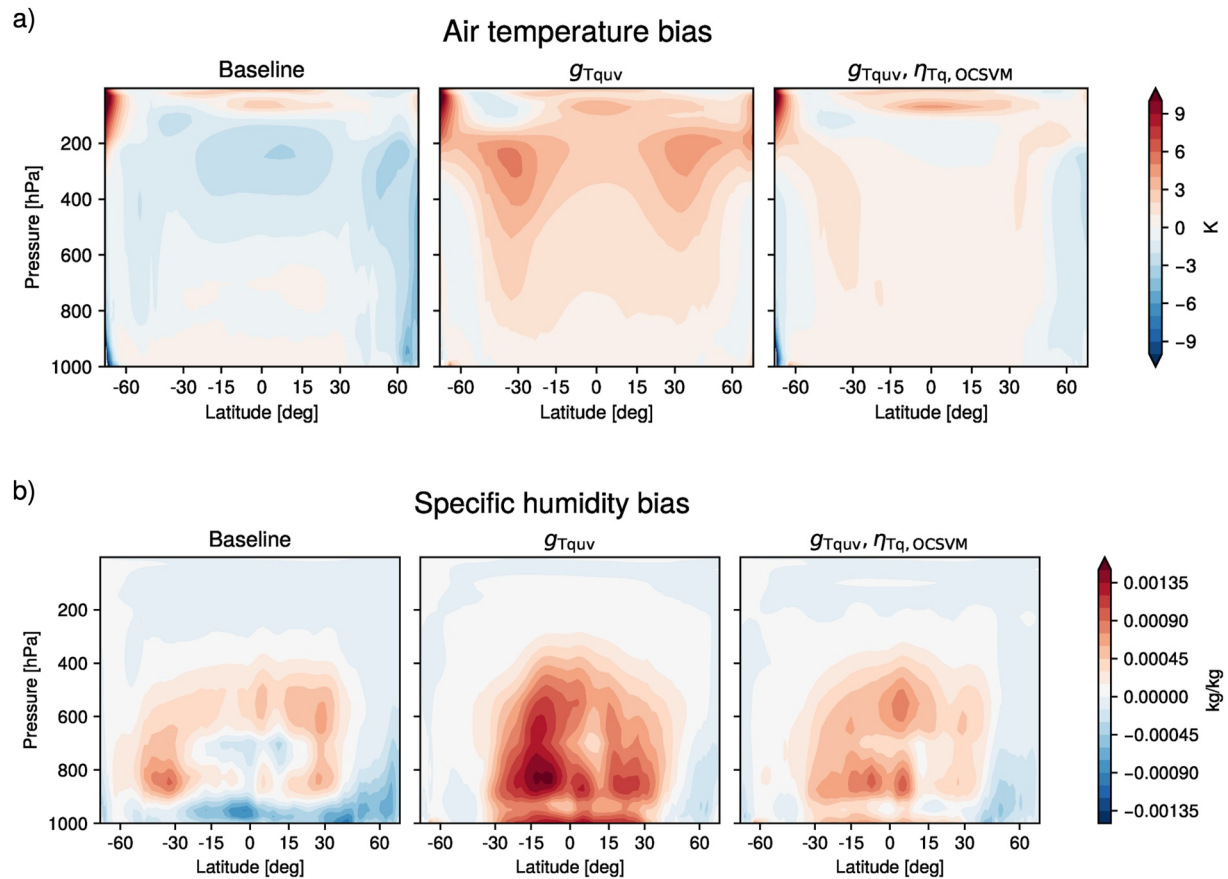


**Figure 3.** Time versus zonal-mean plots visualizing upper-atmospheric temperature biases (against the fine-grid reference simulation) at the 200 hPa pressure level (top) and fractions of novelties identified (bottom) by three different models initialized from random seed 0 (left to right): (1) the ML-corrected climate model  $g_{Tquv}$  without novelty detection, (2)  $g_{Tquv}$  with one-class support vector machine (OCSVM) novelty detection  $\eta_{T, OCSVM}$  using temperature as the input feature, and (3)  $g_{Tquv}$  with OCSVM novelty detection  $\eta_{Tq, OCSVM}$  using temperature and specific humidity as input features.

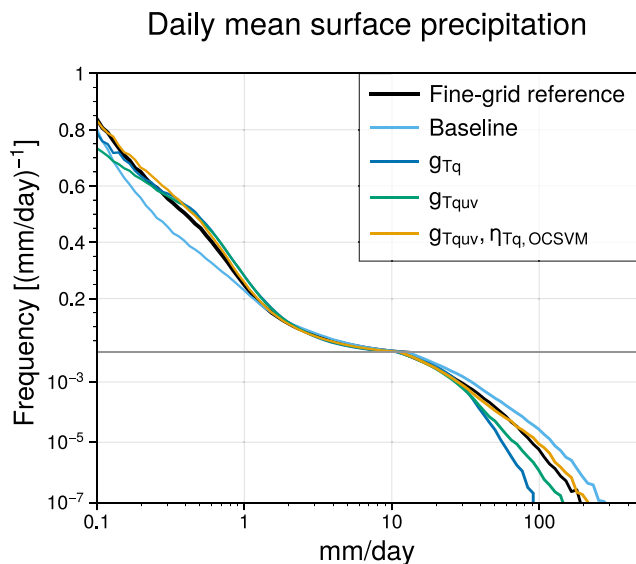
The ML-corrected  $g_{Tquv}$  model without novelty detection develops a significant 5–10 K warm bias in 200 hPa air temperature across latitudes. The temperature-only novelty detection in  $g_{Tquv}, \eta_{T, OCSVM}$  removes the largest magnitude warm bias at the North Pole during boreal summer, but otherwise does not prevent the global warm drift. Though the  $g_{Tquv}, \eta_{T, OCSVM}$  simulation develops 5–10 K biases within the first 16 weeks, the  $\eta_{T, OCSVM}$  activates infrequently as it still classifies these columns' temperature profiles as lying within the training distribution, presumably due to the large weather-associated variability of temperature sampled therein.

The prognostic run in the right column of Figure 3 shows that using specific humidity inputs in addition to temperature inputs is necessary for successful bias reduction via novelty detection. This greatly increases the rate of out-of-sample classification, especially in the tropics. The 200 hPa temperature bias is dramatically reduced out to high latitudes, despite the majority of the novelty detection occurring in the tropics. We speculate that this is due to changes in tropical convection, where the  $\eta_{Tq, OCSVM}$  novelty detector is most active other than extreme polar latitudes.

Figure 4 shows sections of time- and zonal-mean air temperature and specific humidity biases. Instead of the  $g_{Tquv}, \eta_{T, OCSVM}$  run, Figure 4 includes a baseline (no-ML) simulation for comparison, since that is what we are aiming to improve on. The baseline model air temperature is biased low in the tropical stratosphere and throughout the column in high northern latitudes. The ML-corrected  $g_{Tquv}$  model without novelty detection corrects the cold bias at high northern latitudes but develops an overall warm bias that is largest in the extratropical stratosphere.



**Figure 4.** Annual-averaged zonal mean temperature (top) and humidity (bottom) biases plotted over pressure levels, for the baseline model (left) and seed-0  $g_{Tquv}$  models with no novelty detection (center) and with  $\eta_{Tq,OCSVM}$  novelty detection (right).



**Figure 5.** Probability distribution function of daily mean precipitation from all grid columns around the globe, shown for the fine-grid reference, baseline, ML-corrected  $g_{Tquv}$  run without novelty detection, and ML-corrected  $g_{Tquv}, \eta_{Tq,OCSVM}$  run. The y-axis uses linear scaling above 0.01 (mm/day)<sup>-1</sup> and log scaling below.

Adding the  $\eta_{Tq,OCSVM}$  novelty detector on top of the  $g_{Tquv}$  corrections removes most of this stratospheric warm bias.

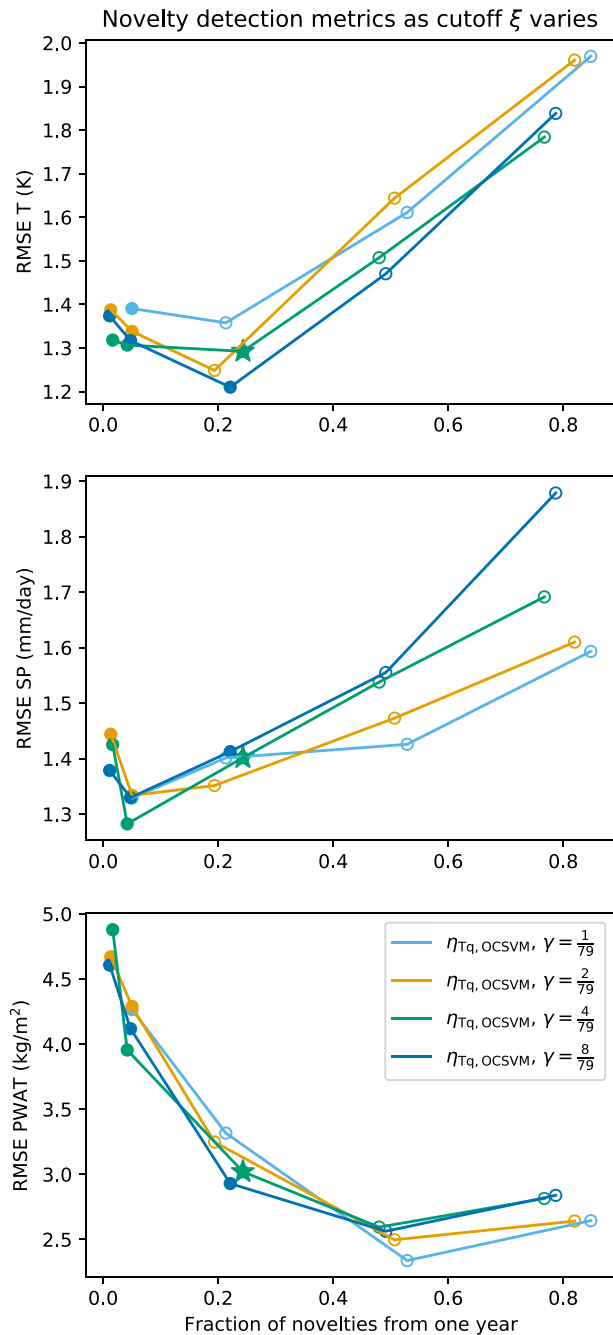
Similarly, the ML-corrected  $g_{Tquv}$  model without novelty detection develops a tropical moist bias in specific humidity that is larger in magnitude than the baseline biases in both the boundary layer and the troposphere. Adding the  $\eta_{Tq,OCSVM}$  novelty detector greatly reduces this bias.

### 3.4. Daily Mean Precipitation Distribution

The ML-corrected  $g_{Tquv}, \eta_{Tq,OCSVM}$  simulation also captures the global-mean probability distribution function (pdf) of daily mean precipitation in the reference fine-grid simulation better than the baseline (no-ML) and  $g_{Tquv}$  approaches (Figure 5). The baseline run underestimates the frequency of low daily mean precipitation below a few mm/day, while the ML-corrected simulations more closely match the fine-grid reference at the low end of the distribution. The baseline run over-estimates the high-precipitation tail of the target pdf, while the  $g_{Tquv}$  run underestimates the pdf in the tail. The  $g_{Tquv}, \eta_{Tq,OCSVM}$  run matches the tail of the global precipitation pdf more closely up to rates over 100 mm/day.

## 4. Varying Novelty Detector Sensitivity

Section 3 considered an OCSVM with  $\gamma = 4/79$  and cutoff  $\xi$  set to the maximum score observed in the training data (See Section 2.3.3 for a more



**Figure 6.** Root mean-square error of time-averaged 850 hPa temperature (top), surface precipitation (center), and precipitable water (bottom) of year-long global C48 simulations, all with ML-correction  $g_{Tquv}$  and novelty detector  $\eta_{Tq, OCSVM}$  with kernel inverse radius  $\gamma \in \left\{ \frac{1}{79}, \frac{2}{79}, \frac{4}{79}, \frac{8}{79} \right\}$ . The plots show each error metric as a function of the total fraction of identified novelties (a monotonic increasing function of  $\xi$ ) on the x-axis. The choices of  $\xi$  are given in Table 1 in Section 2.3.3. The single green star marker represents the one-class support vector machine parameters used in the Results section. Filled markers indicate consistent novelty detectors that classify no more than 5% of the holdout reference training data set as out-of-sample; open circles indicate inconsistent detectors.

thorough discussion of how the value of cutoff  $\xi$  impacts novelty frequency in online simulations.). This model, whether applied only to temperature or to both temperature and humidity, appears to find a consistent “sweet spot” between the baseline run and the ML-corrected run with no novelty detection that reduces the mean-state drifts of both approaches. This section presents a sensitivity study that supports this finding by considering several choices of  $\gamma$  and varying  $\xi$  to adjust the aggressiveness of the novelty detector. We show that these approaches interpolate between the baseline and ML-corrected run as the cutoffs change, and that choosing an intermediate model that categorizes a substantial fraction of samples as novelties balances the trade-off between the quality of temperature and surface precipitation estimates and of precipitable water estimates.

In Figure 6, we consider an ML-corrected model  $g_{Tquv}$  augmented with an OCSVM novelty detector  $\eta_{Tq, OCSVM}$  with various choices of inverse radius parameter  $\gamma$  and cutoff parameter  $\xi$ . We plot the error metrics as a function of the fraction of novelties identified online for each cutoff. We find that an intermediate cutoff balances strong performance on temperature and surface precipitation (for which the ML-correction-only simulation has a lower RMSE than the baseline simulation) and total precipitable water estimates (which are better predicted by the baseline model than the ML-correction-only simulation). Optimal temperature and precipitable water predictions generally occur when approximately 20% and 60% of samples are categorized as novelties, respectively (and hence suppressed). The plots demonstrate that this approach effectively interpolates between those two extreme cases and that the cutoff  $\xi$  used in the preceding section lies near that sweet spot. The figure also demonstrates that different combinations of radius parameter  $\gamma$  and cutoff  $\xi$  result in similarly performing simulations when the fraction of novelties detected is the same.

A potential pitfall of applying novelty detection within simulations is that the detector may falsely flag some columns as novelties and suppress legitimate ML corrections. We try to reduce the occurrence of this behavior by our choice of classification cutoff  $\xi$ . As  $\xi$  is increased, a greater fraction of samples from the reference data set distribution will also be classified as out-of-sample. Detectors are considered inconsistent if they return a significant novelty fraction when evaluated on a holdout set from the training data, as this means that the detector has a higher false positive rate in flagging samples as novelties when they are still within the training distribution. In this analysis we set a false-positive threshold of 5% to determine which  $(\gamma, \xi)$  combinations are consistent OCSVMs. OCSVMs which classify  $>5\%$  of the holdout reference data as out-of-sample are deemed inconsistent and indicated as open circles in Figure 6. These include all detectors that classify less than 75% of online ML-corrected samples as typical and, for certain  $\gamma$ , even detectors classifying up to 95% as such. That is, the best climate performance using this ML correction model is found by using the maximum  $\xi$  consistent with the false-positive threshold on the withheld reference data. The parameters used in the  $\eta_{Tq, OCSVM}$  detector in the preceding sections resulted in a 2.6% false positive rate.

It is also likely that atmospheric states may arise in the coarse-grid model which are consistent with fine-grid model behavior but are wrongly flagged as novelties because the limited time span of the year-long training data set does not fully capture this desired range of behavior. These instances of false positive errors would not be flagged as such by our method described above. They could be reduced by training on longer reference data sets spanning multiple years.



## 5. Conclusion and Future Work

This study demonstrates that applying novelty detection to ML-corrected coarse-grid atmospheric climate models can improve the quality and reliability of their temperature and precipitation estimates. Future efforts using corrective ML models within simulations may find this approach useful for improving forecast skill and avoiding climate drift into states outside the training distribution.

Offline, a novelty detection algorithm trained on samples from a coarsened high-resolution simulation tends to classify more columns as novelties in runs that drift further from the high-resolution reference. When applied online to mask ML-predicted corrective tendencies, the novelty detector maintains or improves the spatial patterns of time-mean surface precipitation rate, lower and upper atmospheric temperature and precipitable water. Furthermore, for an ensemble of ML-corrected simulations (in which each simulation uses an ML model trained with a different random seed initialization of weights), use of novelty detection decreases the spread in model skill across the ensemble. This is a valuable property, since online use of ML parameterizations can be highly sensitive to subtle changes in the offline training, such as random seed (e.g., Wang et al., 2022).

Future work can build on this effort by experimenting with different novelty detection approaches, OCSVM kernels, inputs to  $\eta$ , and methods for integrating the novelty detector into the ML-corrected climate model. Practical implementation of the novelty detector can become a simulation bottleneck if the number of support vectors (2.3.1) is too high. For the settings used in Section 3, the novelty detector roughly doubled the wall clock time per simulation timestep. It would be worth further investigation into how few support vectors are needed to improve ML-corrected simulations online. In addition the more classical ML approaches to novelty detection explored here, future work may consider using neural networks directly for density estimation for the purpose of novelty detection. Finally, further analysis of the character of the out-of-sample behaviors that are being detected by the trained novelty detectors could help us better understand their causes.

## Data Availability Statement

The code used to configure experiments and analyze their results is available at <https://github.com/ai2cm/out-of-sample> (C. Sanford, 2023). The version of the codebase used to train models and run them within coarse-grid simulations is available at <https://github.com/ai2cm/fv3net> (AI2CM, 2023). The coarsened fine-grid data used for initial conditions and in the nudged coarse-grid simulation is available upon request through a Google Cloud Storage “requester pays” bucket.

## Acknowledgments

We thank our reviewers for their valuable feedback, which helped improve the quality of the final paper. This work was started when C. Sanford was a summer intern with the AI2 Climate Modeling group. AI2 is supported by the estate of Paul G. Allen. We appreciate helpful conversations with Daniel Hsu about framing the problem as novelty detection. We thank NOAA-GFDL for running the 1-year X-SHIELD simulation using the Gaea computing system, which was used as the reference fine-grid data set in this work. We also acknowledge NOAA-GFDL, NOAA-EMC, and the UFS community for sharing code, forcing data, and software packages.

## References

- AI2CM. (2023). Fv3net code repository [Software]. Zenodo. <https://doi.org/10.5281/zenodo.7872718>
- Amraee, S., Vafaei, A., Jamshidi, K., & Adibi, P. (2018). Abnormal event detection in crowded scenes using one-class SVM. *Signal, Image and Video Processing*, 12(6), 1115–1123. <https://doi.org/10.1007/s11760-018-1267-z>
- Brenowitz, N. D., Beucler, T., Pritchard, M., & Bretherton, C. S. (2020). Interpreting and stabilizing machine-learning parametrizations of convection. *Journal of the Atmospheric Sciences*, 77(12), 4357–4375. <https://doi.org/10.1175/JAS-D-20-0082.1>
- Brenowitz, N. D., & Bretherton, C. S. (2019). Spatially extended tests of a neural network parametrization trained by coarse-graining. *Journal of Advances in Modeling Earth Systems*, 11(8), 2728–2744. <https://doi.org/10.1029/2019MS001711>
- Brenowitz, N. D., Henn, B., Clark, S., Kwa, A., McGibbon, J., Perkins, W. A., et al. (2020). Machine learning climate model dynamics: Offline versus online performance. In *NeurIPS 2020 workshop on tackling climate change with machine learning*. Retrieved from <https://www.climate-change.ai/papers/neurips2020/50>
- Bretherton, C. S., Henn, B., Kwa, A., Brenowitz, N. D., Watt-Meyer, O., McGibbon, J., et al. (2022). Correcting coarse-grid weather and climate models by machine learning from global storm-resolving simulations. *Journal of Advances in Modeling Earth Systems*, 14(2), e2021MS002794. <https://doi.org/10.1029/2021MS002794>
- Bretherton, C. S., Peters, M. E., & Back, L. E. (2004). Relationships between water vapor path and precipitation over the tropical oceans. *Journal of Climate*, 17(7), 1517–1528. [https://doi.org/10.1175/1520-0442\(2004\)017\(1517:RBWVPA\)2.0.CO;2](https://doi.org/10.1175/1520-0442(2004)017(1517:RBWVPA)2.0.CO;2)
- Breunig, M. M., Kriegel, H.-P., Ng, R. T., & Sander, J. (2000). Lof: Identifying density-based local outliers. In *Proceedings of the 2000 acm sigmod international conference on management of data* (pp. 93–104). Association for Computing Machinery. <https://doi.org/10.1145/342009.335388>
- Chantry, M., Hatfield, S., Dueben, P., Polichtchouk, I., & Palmer, T. (2021). Machine learning emulation of gravity wave drag in numerical weather forecasting. *Journal of Advances in Modeling Earth Systems*, 13(7), e2021MS002477. <https://doi.org/10.1029/2021ms002477>
- Chen, T.-C., Penny, S. G., Whitaker, J. S., Frolov, S., Pincus, R., & Tulich, S. (2022). Correcting systematic and state-dependent errors in the NOAA fv3-gfs using neural networks. *Journal of Advances in Modeling Earth Systems*, 14(11), e2022MS003309. <https://doi.org/10.1029/2022MS003309>
- Cheng, K.-Y., Harris, L., Bretherton, C., Merlis, T. M., Bolot, M., Zhou, L., et al. (2022). Impact of warmer sea surface temperature on the global pattern of intense convection: Insights from a global storm resolving model. *Geophysical Research Letters*, 49(16), e2022GL099796. <https://doi.org/10.1029/2022GL099796>
- Clark, S. K., Brenowitz, N. D., Henn, B., Kwa, A., McGibbon, J., Perkins, W. A., et al. (2022). Correcting a coarse-grid climate model in multiple climates by machine learning from global 25-km resolution simulations. *Earth and Space Science Open Archive*, 46. <https://doi.org/10.1002/essoar.10511517.1>

- Hodge, V., & Austin, J. (2004). A survey of outlier detection methodologies. *Artificial Intelligence Review* (Vol. 85–126). Retrieved from <https://eprints.whiterose.ac.uk/767/>
- IPCC. (2021). Summary for policymakers [book section]. In V. Masson-Delmotte, P. Zhai, A. Pirani, S. L. Connors, C. Pean, S. Berger, et al. (Eds.), *Climate change 2021: The physical science basis. Contribution of working group I to the sixth assessment report of the intergovernmental panel on climate change* (pp. 3–32). Cambridge University Press. <https://doi.org/10.1017/9781009157896.001>
- Krasnopolsky, V. M., Fox-Rabinovitz, M. S., & Belochitski, A. A. (2010). Development of neural network convection parameterizations for numerical climate and weather prediction models using cloud resolving model simulations. In *The 2010 international joint conference on neural networks* (pp. 1–8). IJCNN. <https://doi.org/10.1109/IJCNN.2010.5596766>
- Krasnopolsky, V. M., Fox-Rabinovitz, M. S., Tolman, H. L., & Belochitski, A. A. (2008). Neural network approach for robust and fast calculation of physical processes in numerical environmental models: Compound parameterization with a quality control of larger errors (Vol. 21). <https://doi.org/10.1016/j.neunet.2007.12.019>
- Kwa, A., Clark, S. K., Henn, B., Brenowitz, N. D., McGibbon, J., Watt-Meyer, O., et al. (2022). Machine-learned climate model corrections from a global storm-resolving model: Performance across the annual cycle. <https://doi.org/10.1002/essoar.10512393.1>
- Nairac, A., Townsend, N., Carr, R., King, S., Cowley, P., & Tarassenko, L. (1999). A system for the analysis of jet engine vibration data. *Integrated Computer-Aided Engineering*, 6(1), 53–66. <https://doi.org/10.3233/ica-1999-6106>
- Putman, W. M., & Lin, S.-J. (2007). Finite-volume transport on various cubed-sphere grids. *Journal of Computational Physics*, 227(1), 55–78. <https://doi.org/10.1016/j.jcp.2007.07.022>
- Rasp, S., Pritchard, M. S., & Gentine, P. (2018). Deep learning to represent subgrid processes in climate models. *Proceedings of the National Academy of Sciences*, 115(39), 9684–9689. <https://doi.org/10.1073/pnas.1810286115>
- Sanford, C. (2023). Out-of-sample experiment configuration and analysis code repository [Software]. Zenodo. <https://doi.org/10.5281/zenodo.7872723>
- Sanford, C. H., Kwa, A., Watt-Meyer, O., Clark, S., Brenowitz, N., McGibbon, J., & Bretherton, C. (2022). Improving the predictions of ml-corrected climate models with novelty detection. In *Neurips 2022 workshop on tackling climate change with machine learning*. Retrieved from <https://www.climatechange.ai/papers/neurips2022/10>
- Schölkopf, B., Platt, J., Shawe-Taylor, J., Smola, A., & Williamson, R. (2001). Estimating support of a high-dimensional distribution. *Neural Computation*, 13(7), 1443–1471. <https://doi.org/10.1162/089976601750264965>
- Schölkopf, B., Smola, A., Williamson, R. C., & Bartlett, P. L. (2000). New support vector algorithms. *Neural Computation*, 12(5), 1207–1245. <https://doi.org/10.1162/089976600300015565>
- Sommer, C., Hoefler, R., Samwer, M., & Gerlich, D. W. (2017). A deep learning and novelty detection framework for rapid phenotyping in high-content screening. *Molecular Biology of the Cell*, 28(23), 3428–3436. PMID: 28954863. <https://doi.org/10.1091/mbc.e17-05-0333>
- Song, H.-J., Roh, S., & Park, H. (2021). Compound parameterization to improve the accuracy of radiation emulator in a numerical weather prediction model. *Geophysical Research Letters*, 48(20), e95043. <https://doi.org/10.1029/2021gl095043>
- Tan, Y., Niu, C., Tian, H., Hou, L., & Zhang, J. (2019). A one-class svm based approach for condition-based maintenance of a naval propulsion plant with limited labeled data. *Ocean Engineering*, 193, 106592. <https://doi.org/10.1016/j.oceaneng.2019.106592>
- Van Aelst, S., & Rousseeuw, P. (2009). *Minimum volume ellipsoid* (Vol. 1, pp. 71–82). Wiley Interdisciplinary Reviews: Computational Statistics. <https://doi.org/10.1002/wics.19>
- Wang, X., Han, Y., Xue, W., Yang, G., & Zhang, G. J. (2022). Stable climate simulations using a realistic general circulation model with neural network parameterizations for atmospheric moist physics and radiation processes. *Geoscientific Model Development*, 15(9), 3923–3940. <https://doi.org/10.5194/gmd-15-3923-2022>
- Watt-Meyer, O., Brenowitz, N. D., Clark, S. K., Henn, B., Kwa, A., McGibbon, J., et al. (2021). Correcting weather and climate models by machine learning nudged historical simulations. *Geophysical Research Letters*, 48(15), e2021GL092555. <https://doi.org/10.1029/2021GL092555>
- Woelfle, M. D., Yu, S., Bretherton, C. S., & Pritchard, M. S. (2018). Sensitivity of coupled tropical pacific model biases to convective parameterization in cesm1. *Journal of Advances in Modeling Earth Systems*, 10(1), 126–144. <https://doi.org/10.1002/2017MS001176>
- Yang, K., Kpotufe, S., & Feamster, N. (2021). An efficient one-class SVM for anomaly detection in the internet of things. CoRR, abs/2104.11146. Retrieved from <https://arxiv.org/abs/2104.11146>
- Yuval, J., & O’Gorman, P. A. (2020). Stable machine-learning parameterization of subgrid processes for climate modeling at a range of resolutions. *Nature Communications*, 11(1), 3295. <https://doi.org/10.1038/s41467-020-17142-3>
- Zhang, G. J., & Wang, H. (2006). Toward mitigating the double ITCZ problem in NCAR CCSM3. *Geophysical Research Letters*, 33(6), L06709. <https://doi.org/10.1029/2005GL025229>
- Zhou, L., Lin, S.-J., Chen, J.-H., Harris, L. M., Chen, X., & Rees, S. L. (2019). Toward convective-scale prediction within the next generation global prediction system. *Bulletin of the American Meteorological Society*, 100(7), 1225–1243. <https://doi.org/10.1175/BAMS-D-17-0246.1>