# Multi-Year-to-Decadal Temperature Prediction using a Machine Learning Model-Analog Framework

M. A. Fernandez[1] and Elizabeth A. Barnes[1]

[1]Department of Atmospheric Science, Colorado State University, Fort Collins, CO, USA

**Key Points:**

- We use machine learning to identify important precursor regions for matching in a model-analog prediction framework.
- We predict 2-meter temperature with lead times from one to ten years and for five different regions.
- We find improved performance over other analog prediction methods and over initialized decadal predictions.

Corresponding author: M. A. Fernandez, `mafern@colostate.edu`

**Abstract**

Multi-year-to-decadal climate prediction is a key tool in understanding the range of potential regional and global climate futures. Here, we present a framework that combines machine learning and analog forecasting for predictions on these timescales. A neural network is used to learn a mask, specific to a region and lead time, with global weights based on relative importance as precursors to the evolution of that prediction target. A library of mask-weighted model states, or potential analogs, are then compared to a single mask-weighted observational state. The known future of the best matching potential analogs serve as the prediction for the future of the observational state. We match and predict 2-meter temperature using the Berkeley Earth Surface Temperature dataset for observations, and a set of CMIP6 models as the analog library. We find improved performance over traditional analog methods and initialized decadal predictions.

## Plain Language Summary

Accurate prediction of the climate, from the next year to the next decade, is important for adaptation and mitigation in many societal sectors. In particular, temperature plays a key role on these time scales for agriculture, human health, and infrastructure planning. In this work, we develop and test a framework for making predictions on multi-year-to-decadal timescales, utilizing machine learning and an established method called analog forecasting. Analog forecasting compares the climate state from which one wants to predict to a library of potential analog states. The forecast is then the known future of the best matching analog states, which presumes that the analogs evolve similarly to the real world. Here, our analogs come from a subset of Earth system model runs from the CMIP6 dataset. The machine learning component contributes by learning the regions that are most important in determining if two climate states will evolve similarly, ensuring that the analog selection encodes information specific to the prediction task. We find that our method improves on current prediction methods for these timescales in several regions and lead times.

## 1 Introduction

Regional climate prediction from years to decades is of growing interest to the climate science community, governing bodies around the world, and the general population. In particular, skillful predictions of the climate one to ten years in the future are necessary for understanding both the future climatological mean and potential extremes (Hermanson et al., 2022; IPCC, 2023a). These predictions are therefore vital for mitigation and adaptation planning on regional scales (Khasnis & Nettleman, 2005; Solaraju-Murali et al., 2022; Dunstone et al., 2022; IPCC, 2023c, 2023b).

There are many methods aimed at understanding the future climate on multi-year to decadal timescales. Many of these make use of the large collection of available Earth System Models (ESMs) that have been developed at research centers around the world. These include the Shared Socioeconomic Pathways (SSP) projections (Eyring et al., 2016), which explore a range of potential climate futures out to the year 2100. However, SSPs were not designed nor intended to be used as predictions.

Initialized ESMs (IESMs) use observations to initialize physics-based dynamical models, then run these simulations, often out to ten years (Meehl et al., 2021). If the model captures the underlying physics of the Earth system, then this initialization should lead to skillful realizations. However, due to model biases and unresolved physics, IESMs suffer from climate drift, where the simulations converge back to the model mean climate (Meehl et al., 2021, 2022). IESMs can also be costly to run, making large ensembles difficult and reducing their effectiveness for the prediction of extremes.

Work in the last several years (Befort et al., 2020; Mahmood et al., 2021; Befort et al., 2022; De Luca et al., 2023; Donat et al., 2024) has used both observations and IESMs to constrain climate projections (SSPs) of the next 10 to 40 years. These studies match the projections to either observations or IESMs over a common period for a specified region, keeping only the best matching members to create a constrained sub-ensemble. These constrained sub-ensembles have been shown to have better skill than the full ensemble for up to 20 years.

An initialization method that has a long history is analog forecasting (Lorenz, 1969). In analog forecasting, matching is based on a single state (model, member, and time), and the known future of the best matching state or states is the prediction. As opposed to other initialized prediction methods, analog forecasting is time-agnostic, e.g., the prediction for 2025 does not need to come from a projection of 2025. Initially, analog forecasting was used for weather prediction using historical observations as the analogs for the current state of the system (Bergen & Harnack, 1982).

There have been several recent studies exploring the use of models in an analog forecasting framework, called model-analogs (Ding et al., 2018; Rader & Barnes, 2023; Ding & Alexander, 2023; Toride et al., 2024; Acosta Navarro et al., 2025). The underlying premise of model-analogs is that there already exists a large catalog of ESMs that can be used to match to observations, and from which analog predictions can be made. Figure 1 shows a time series of obervations, a suite of models covering the same time period, and a posthoc selected set of analogs. The posthoc selection shown in Figure 1 is useless in a forecasting sense, but illustrates that within currently available model runs, states exist that closely match observations not just for a single year, but for multiple years.
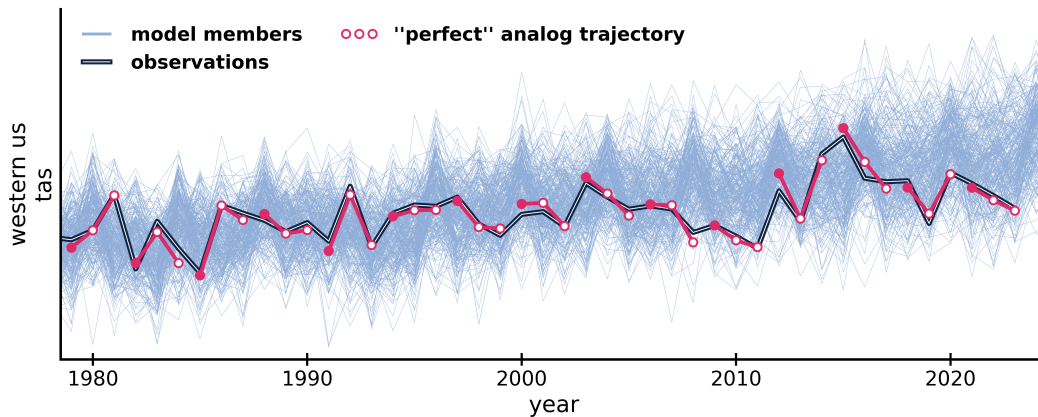


**Figure 1.** Posthoc selection of best matching analog to observations in three year chunks. This is not usable as a prediction.

Here, we build on work from Rader and Barnes (2023), which used machine learning to identify the most important precursor regions for matching in a model-analog framework. We use a similar machine learning set up, along with a new matching criteria, and an expanded dataset that includes members from 29 Coupled Model Intercomparison Project Phase 6 (CMIP6; Eyring et al., 2016) models. We make predictions for lead times from one to ten years in multiple regions. We evaluate these predictions against alternative analog methods, as well as bias-corrected IESMs, using multiple metrics to capture performance on the mean trend, variability, and pattern.

## 2 Prediction Framework

Our aim is to capitalize on existing Earth system model (ESM) simulations for prediction on multi-year-to-decadal timescales. We do this using an analog forecasting framework. Here, the library of analogs is composed of model simulations from the CMIP6 suite, and the matching method uses a learned mask of weights that emphasizes important precursor regions for the specific prediction task (target region, lead time, and variable). We match on, and predict, annual mean 2-meter temperature. Some example masks for different regions for a lead time of five years are shown in Figure 2.
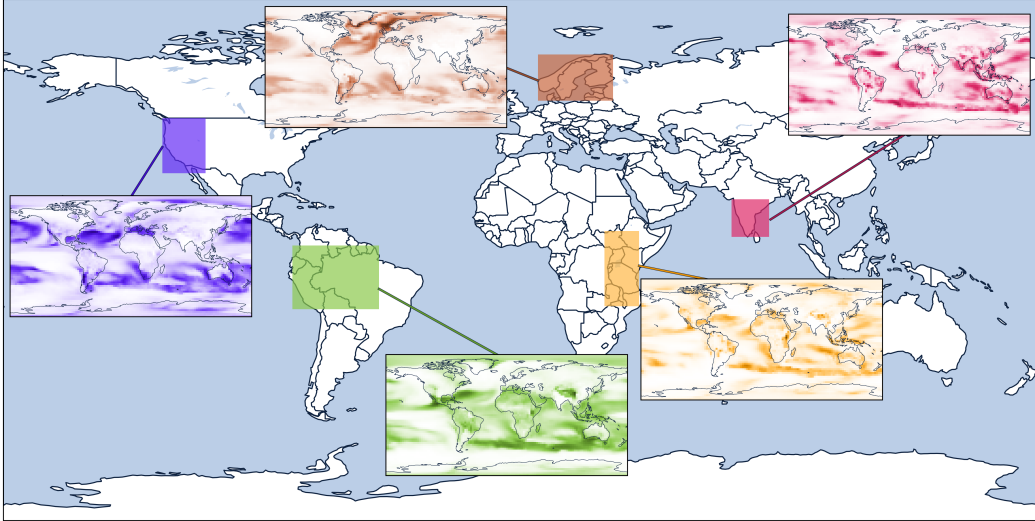


**Figure 2.** Model learned masks for annual mean 2-meter temperature for five regions: the western United States, the Amazon, northern Europe, the great lakes of Africa, and southern India. Each region (in fact, each grid cell) has a unique mask highlighting global locations that are important markers for the evolution of that region, on that time scale.

### 2.1 Data

Our dataset includes 2-meter temperature from 29 ESMs from the CMIP6 suite of simulations, detailed in Table 1. All model runs are obtained from the Earth System Grid Federation (ESGF; Cinquini et al., 2014). All members include the historical period from 1850-2014, while a subset contain the Shared Socioeconomic Pathways (SSP) projections corresponding to SSP-3.7.0 (usually 2015-2100, but some models do not reach 2100). In total, there are 285 simulations that include both the historical period and SSP projections, with an additional 183 simulations that include only the historical period. Observations of surface temperature are from the Berkeley Earth Surface Temperature (BEST) dataset (Rohde & Hausfather, 2020).

All data is resampled to annual averages, leaving over $100,000$ states in the analog library. Data is regridded to the coarsest native resolution among the models, 2.77 (latitude) by 2.81 (longitude) degrees, which translates to 65 latitudes and 128 longitudes. Higher resolutions (1.5 and 2 degrees) were tested for five year predictions of the western United States, but did not improve performance. Data is normalized by baselining on a 30 year $(1961-1990)$ mean for each member, which allows even warm biased models to be useful members of the analog library, while maintaining the mean climate trend.

**Table 1.** CMIP6 dataset: total number of members for each model, as well as number (in parenthesis) that include their associated SSP-3.7.0 projections.

| model | members | model | members | model | members |
|---|---|---|---|---|---|
| ACCESS-CM2 | 10 (5) | ACCESS-ESM1-5 | 40 (40) | AWI-CM-1-1 | 5 (5) |
| AWI-ESM-1-1 | 1 (0) | BCC-CSM2 | 3 (1) | BCC-ESM1 | 3 (0) |
| CAMS-CSM1 | 3 (2) | CESM2 | 10 (3) | CESM2-WACCM | 3 (3) |
| CMCC-CM2-HR4 | 1 (0) | CMCC-CM2-SR5 | 11 (1) | CMCC-ESM2 | 1 (1) |
| CNRM-CM6-1 | 24 (6) | CanESM5 | 60 (45) | CanESM5-1 | 72 (20) |
| CanESM5-CanOE | 3 (3) | GISS-E2-1-G | 47 (23) | GISS-E2-2-G | 11 (5) |
| MIROC6 | 50 (50) | MIROC-ES2H | 3 (0) | MIROC-ES2L | 30 (10) |
| MPI-ESM1-2-HR | 10 (10) | MPI-ESM1-2-LR | 30 (30) | MRI-ESM2-0 | 11 (5) |
| NorESM2-LM | 3 (3) | NorESM2-MM | 2 (1) | TaiESM | 2 (1) |
| UKESM1-0-LL | 18 (11) | UKESM1-1-LL | 1 (1) | | |

**Table 2.** Initialized Earth System Models (IESMs): Both models use full field initialization beginning in November of each year and are run out to ten years (Meehl et al., 2021).

| model | members | initialized components |
|---|---|---|
| CESM1-1-CAM5-CMIP5 | 40 | Ocean |
| CMCC-CM2-SR5 | 10 | All |

To compare our predictions with established methods, we obtain two IESMs, described in Table 2. The IESMs are bias corrected by adding a lead time dependent bias to each member. The bias is the mean difference between observations and each ensemble member for the 15 year period preceding and including the initialization year. There are many methods for correcting model drift in IESMs (Meehl et al., 2022). While more complicated correction methods may improve the IESMs performance, we opt for this correction for two reasons. First, it is a simple method of correcting the IESMs. Second, for clarity we present our method without making similar corrections to the analog library. To avoid unduly favoring the IESMs, we maintain the simple IESM correction. IESMs are then processed using the same steps as the analog library.

## 2.2 Mask of Weights

In order to identify precursor regions on which to match an initial state with a library of analogs, we turn to machine learning. The CMIP6 data is split into three sets: an analog library, a training library, and a validation library. The analog library is split off first, so that it contains at least one member from every model (there are no repeated members between the three libraries). For each model, these libraries contain at most five members for the analog library, three members for the training library, and two members for the validation library, depending on availability. Each of these libraries is processed into an input set and a target set. The target region is selected for the target data. Both the input and target data are normalized by their mean and standard deviation, and shifted by the prediction lead time.

Figure 3 illustrates the machine learning task we employ to learn these precursor regions. Two random states (global maps, one from the analog library, one from the training library) are selected and their difference is multiplied by a mask of weights (initialized as all ones). For training, the batch size used is 64 and the activation function for the mask is relu. The weighted mean-squared error (MSE) of these two states is then
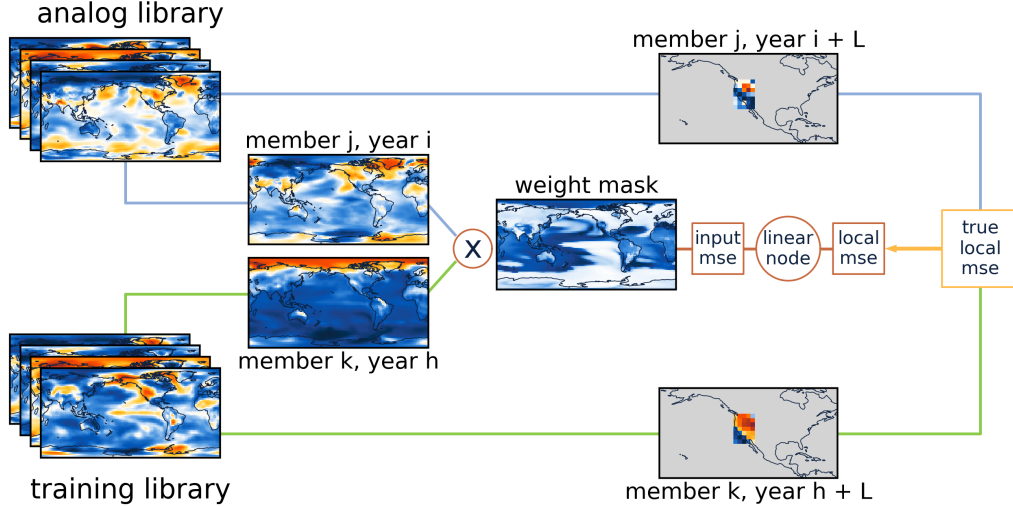
**Figure 3.** Schematic illustrating the machine learning component of the prediction framework. The learned mask identifies global precursors for predicting a specific variable for a given lead time (L in the figure) for a specific region (western United States in the figure).

passed through a single linear node which predicts the corresponding future target region MSE. The mask of weights is updated through backpropagation (learning rate of 0.001) until a strong correspondence between weighted input MSE and future target region MSE is established. We employ early stopping (patience of 50 epochs and minimum change of 0.0005) using the MSE computed between states selected from the analog library and validation library, which contains 2500 samples. The loss is the error between the predicted and true future target MSE.

In addition to the model learned mask, we train a second mask using transfer learning (Pan & Yang, 2010; Ghani et al., 2024). We update the model learned mask by retraining on models (training and validation library) and BEST observations (analog library, includes years up to 2008). The learning rate is reduced to 0.0001, and a trainable dense layer with five nodes and elu activation are added to the architecture between the input MSE and linear node. We manually stop the training after 150 epochs, to avoid the mask forgetting too much of the model learned mask weights.

### 2.3 Analog Selection

From this point on there is no machine learning, only the learned mask is retained for the remainder of the prediction framework. The full CMIP6 data set is preprocessed using the same steps outlined in Section 2.1, without splitting into multiple libraries.

Analog selection is usually done by matching the single initialization state to all of the individual potential analog states. This means finding the lowest MSE between the initialization state and the library of potential analogs. The match is often calculated either everywhere on the globe (a global mask) or in the region of interest (a regional mask). We compare the performance of our method to these two baselines in Section 3. In contrast, we use our learned mask and we match over multiple years (which we call "tethering") to find the best analogs.

Figure 4 is a schematic of our analog selection method. In the case of tethering, the initialization state has a weight of one, while previous years are inverse weighted, e.g.,
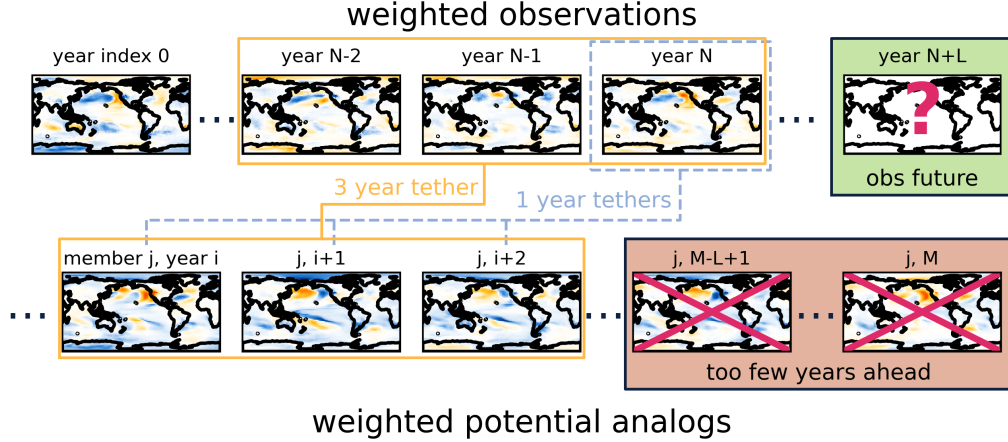
**Figure 4.** Schematic illustrating the method we use to select analogs. Rather than match only the initialization state to the analog library (blue dashed lines), we enforce that years prior to the initialization also match (yellow solid lines, an example three year match). We call this tethering and it helps avoid selecting by-chance matches that may have very different evolutions.

1/2, 1/3 in the example shown. Throughout, we use a tether length of two years, though the length of the optimal tether will likely be region and lead time specific, and could be determined using a left out model. We leave this for future work. An illustrative example of the year of the selected analogs versus each prediction year is shown in Figure A8 for the model learned mask and a regional mask.

### 2.4 Evaluation Metrics

In order to comprehensively evaluate the performance of our prediction framework, we use several metrics. As a base metric, we calculate the MSE between the prediction and truth over the target region. We also evaluate our method using the continuous ranked probability score (CRPS; Gneiting & Raftery, 2007). The CRPS measures how representative a predicted distribution is, given a true value, allowing us to quantify the performance of the analog distribution. CRPS collapses to mean absolute error in the case of a deterministic prediction. Lower values of CRPS indicate better performance.

Additional metrics shown in the Supplementary Material are the Earth Mover's Distance (EMD, calculated using the pyemd[1] package, Pele & Werman, 2008, 2009) and a measure of class accuracy. The EMD measures similarity between two distributions (Rubner et al., 1998); in our case, the predicted and true spatial temperature distribution at each time. The EMD is an optimal transport problem, where, in our case, predicted temperatures are shifted around the target region until the predicted map matches the true map. Smaller values of EMD indicate the prediction was closer to the truth. In addition to redistributing temperatures to match, some excess temperature (bias) may need to be added or discarded, and a penalty can be added for this case. We use a penalty of one, meaning for every degree added or removed, one is added to the final EMD value. The EMD allows us to quantify how well our analogs are matching the temperature patterns in the region of interest.

---

[1] https://pypi.org/project/pyemd/

For the class accuracy, we coarsen the truth and the prediction into a number of bins (bin edges are spatially and temporally determined), then calculate what fraction of the predictions fall in the same bin as the truth. In the Supplementary Material, we show results for four classes. The class accuracy allows us to quantify how well the analogs capture temperature more broadly.

## 3 Results

Using the MSE and CRPS metrics (Section 2.4), we now compare the learned mask analogs to global and regional mask analogs, as well as IESMs. Results for additional metrics (EMD, class accuracy) can be seen in Supplementary Material Figures A5-A7.

Throughout, predictions are for a single annual average of the target region map, e.g., the 5 year prediction initialized in 2010 is for the year 2015 only. We treat each prediction as continuous, e.g., the best matching analog for each year is combined into a single time series map prediction, as is the second best match for each year, and so on. We refer to these time series map predictions as analog forecasts, which we use individually (for metric distributions) or to create an ensemble mean prediction. Example predictions, averaged over the target region, are shown in Figure 5. Individual analog forecasts allow us to assess predictions that have similar variability to observations, while the analog ensemble mean prediction helps determine how well we are capturing average trends (Kim et al., 2025). All summary metrics shown are the mean metric over the target region and time. We first assess performance for lead times 1 to 10 years for a single region (the western United States). We then assess performance for a single lead time (5 years) in five different regions (shown in Figure 2).
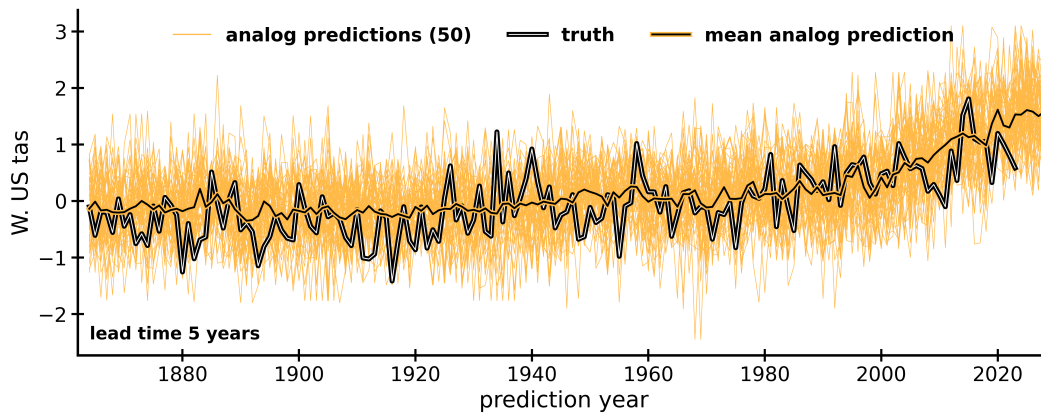


**Figure 5.** Example region-averaged analog predictions and observations for the western United States for a lead time of 5 years. The 50 analog ensemble mean is shown (yellow-black), along with each of the 50 individual analog forecasts (thin, yellow), and the truth from the BEST observational dataset (white-black).

### 3.1 Western United States 1-10 Year Predictions

The learned masks for western United States 2-meter temperature prediction are lead time dependent. We show three examples of the model learned and transfer learned masks on the left and right sides of Figure 6, respectively. Focusing first on the model learned masks, we see that the region itself is usually an important precursor region. Though varying in the specific pattern, the mid-latitudes are consistently important precursor regions, while the tropics are less important. The Southern Ocean is not an important
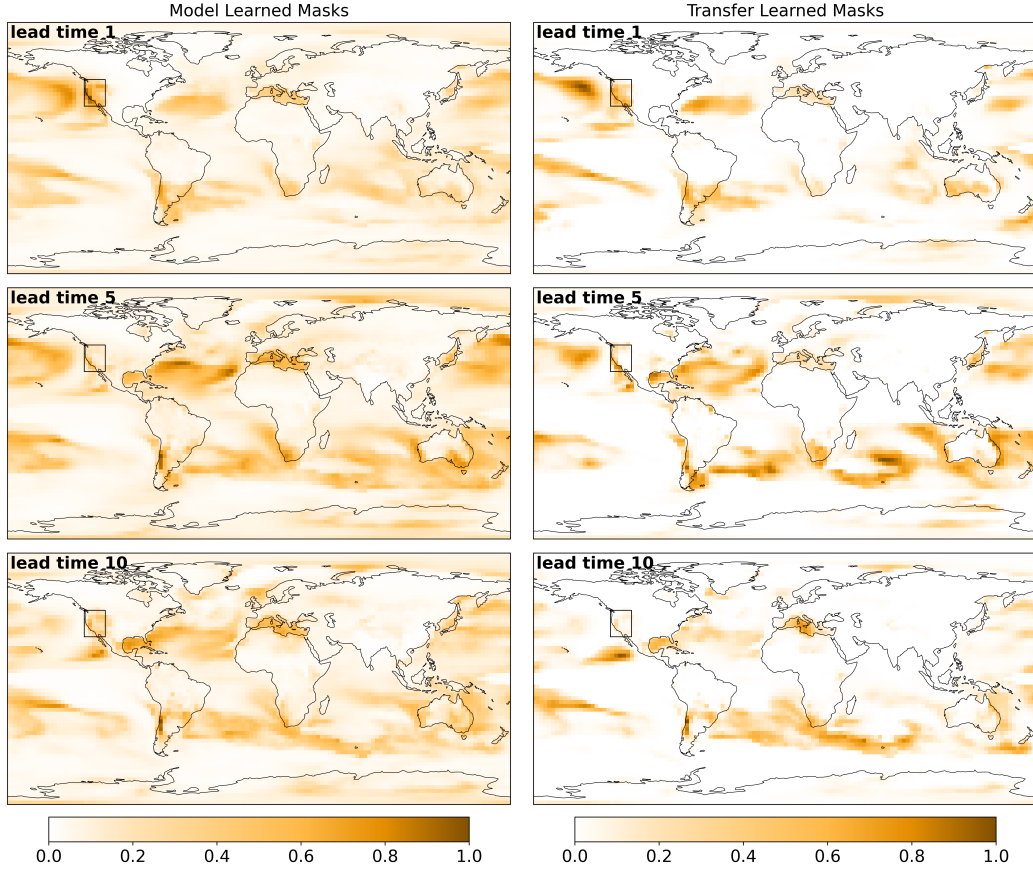
**Figure 6.** Weighted masks for the western United States region (the black bounded area), scaled by the maximum weight to highlight the pattern. The left panels show model learned masks for lead times of 1, 5, and 10 years. The right panels show transfer learned masks (see Section 2.2) for the same lead times.

precursor region for any of the lead times, while the Mediterranean Sea is important for all lead times. The Gulf of Mexico is important for lead times greater than 1 year. North America, outside the western United States region, Asia, Europe, and northern Africa all have low importance. Antarctica often has some importance, while the Arctic Ocean has negligable weights.

The strongest exception to these general trends occur for lead times of 4 and 8 years (not shown here, masks for all lead times can be seen in Figures A1 and A2 in the Supplementary Material). For those lead times, North America, Asia, Europe, and northern Africa have more importance. Antarctica has nearly zero weights, while the Arctic Ocean is more important than at the other lead times.

The transfer learned masks in Figure 6 show the masks after finetuning on a subset of the BEST dataset. The transfer learned masks have primarily been refined, with important regions made even more important, and low importance regions reduced further, indicating that the observations prefer more sparsity.

Figure 7 shows the CRPS and MSE (mean over the time period 1864-2023) for analog forecasts of the western United States, with lead times 1 to 10 years. We use 50 analog forecasts throughout. MSE and CRPS versus the number of analogs used is shown
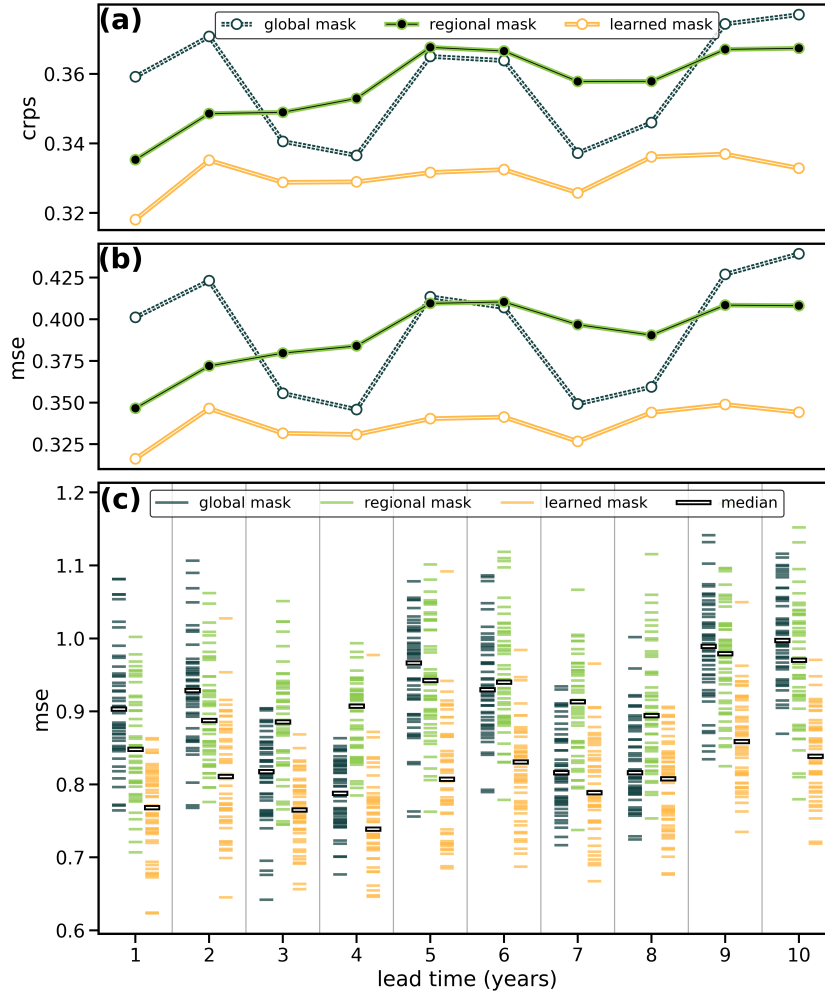
**Figure 7.** Western United States prediction metrics for lead times 1 to 10 years, covering the period 1864-2023. (a) CRPS for 50 analog forecast distributions. (b) MSE for the 50 analog ensemble mean predictions. (c) MSE for each distribution of 50 analog forecasts, as well as the median MSE from these distributions.

in Supplementary Material Figures A3 and A4. For years before 1956, the BEST dataset does not have data at every global location, thus the learned and global analogs are matched on only the available locations for those years.

The metrics in Figure 7 are all for analog forecasts, but the analogs are selected by matching either everywhere (global), in just the western United States (regional), or using our model learned mask of weights. At all lead times, the model learned mask produces a better distribution of analogs according to the CRPS metric, shown in panel (a). Likewise, the learned mask produces an ensemble mean prediction that better tracks the true mean regional trend at all lead times, according the MSE of the ensemble mean prediction, shown in panel (b). Finally, the analog forecast MSE distribution, shown in panel (c), using the model learned mask is lower than the other masks, indicating that the learned masks selects analogs that are better representations of the regional variability.

In Figure 8, we compare the performance of our model learned mask, along with our transfer learned mask, to a set of 50 IESM members. The time period covered for
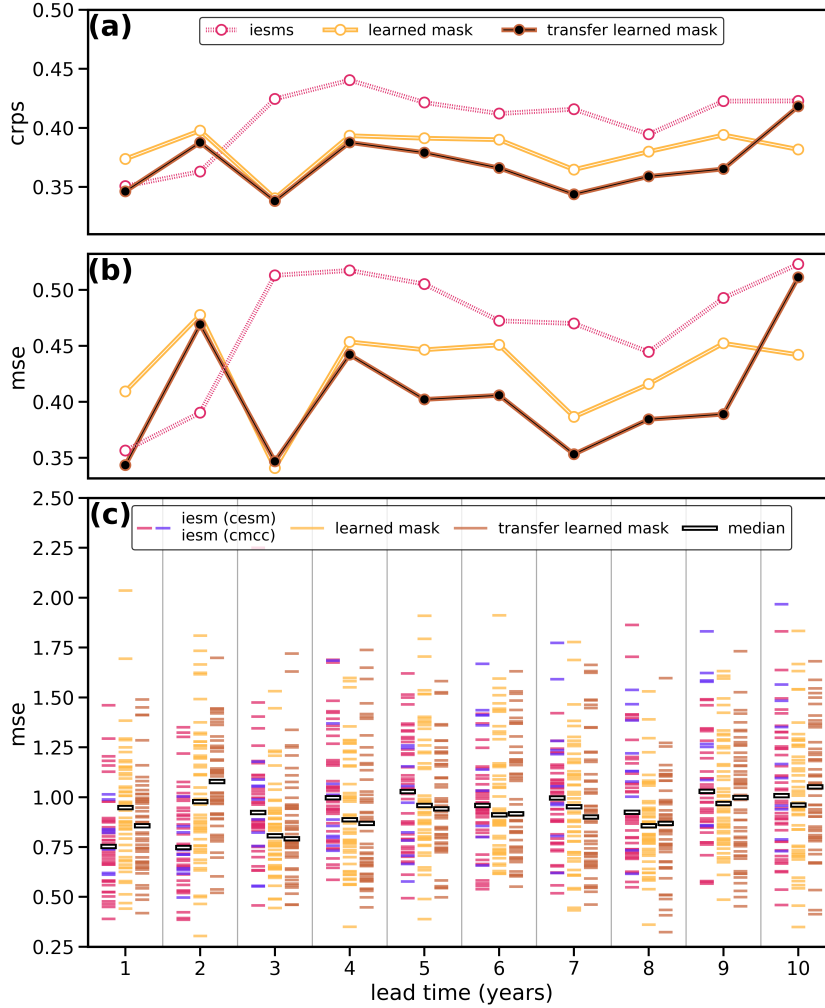
**Figure 8.** Same as Figure 7, but comparing the model learned mask, the transfer learned mask, and the IESMs. The period covered is 2009-2018, which allows observations up to 2008 to be used for the transfer learning.

these metrics is more limited for two reasons. For the transfer learning, we need as many samples as we can get while still leaving some years for testing. We use BEST observations up to 2008. For the IESMs, the latest common initialization year is 2017, so 2018 is the final year for which we have predictions for every lead time. Thus, the period examined in Figure 8 is ten years, 2009-2018.

For both CRPS and MSE, the IESMs outperform our learned masks for lead times 1 and 2 years. For lead times longer than 2 years, our learned masks outperform the IESMs. The transfer learned mask improves upon the model learned mask in several cases. For both the CRPS and ensemble mean MSE (panels (a) and (b)), it is as good or better than the model learned mask for all lead times except 10. For the analog forecast MSE distribution (panel (c)), the transfer learned mask improves on the model learned mask for lead times 1, 3, 4, 5, and 7. Our focus here is not on transfer learning, thus we have not optimized the transfer learning process (number of observation samples, training epochs, updated architecture). Based on these results, incorporating observations in the analog framework is a promising direction for future work.

### 3.2 Regional 5 Year Predictions

We now explore the performance of our method in different regions. To the western United States we add four additional regions: the Amazon, northern Europe, southern India, and the African Great Lakes. All predictions in this section are for a lead time of 5 years, and we explore the model learned mask only, dropping the transfer learned mask so that results can be evaluated over a longer time period.
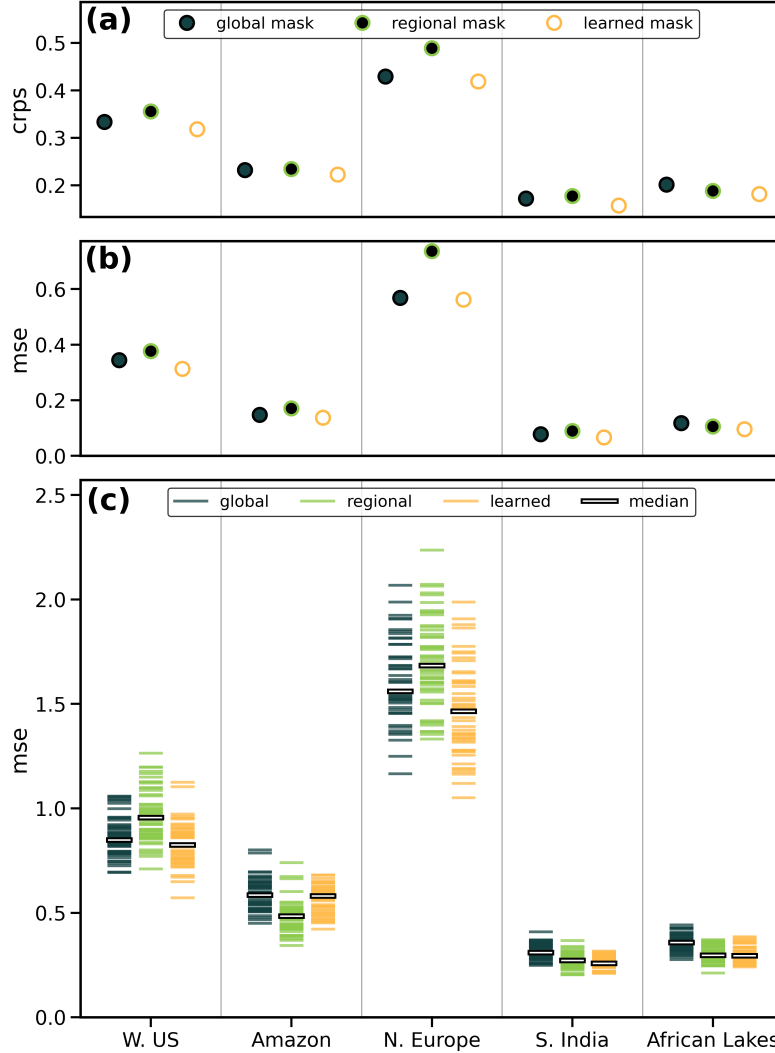


**Figure 9.** 5 year prediction metrics for five different regions (see Figure 2) covering the period 1956-2023. (a) CRPS for 50 analog forecast distributions. (b) MSE for the 50 analog ensemble mean predictions. (c) MSE for each distribution of 50 analog forecasts, as well as the median MSE from these distributions.

We once again start by comparing to alternative analog methods in Figure 9, which shows the same metrics as the figures in Section 3.1. For the CRPS and ensemble mean MSE (panels (a) and (b)), the learned mask outperforms the global and regional mask in all of the regions, with the global mask generally outperforming the regional. For the analog forecast MSE distribution (panel (c)), the learned mask outperforms the global and regional mask in four out of the five regions. The Amazon region is the exception,

where the regional mask performs best. This indicates that local conditions in the Amazon region are very important for predicting future variability of that region. The learned mask (Figure 2) does place a lot of weight in the Amazon region, but highlights other regions as more important. This may mean the learned masks could be improved with the addition of a loss term that more directly accounts for variability.
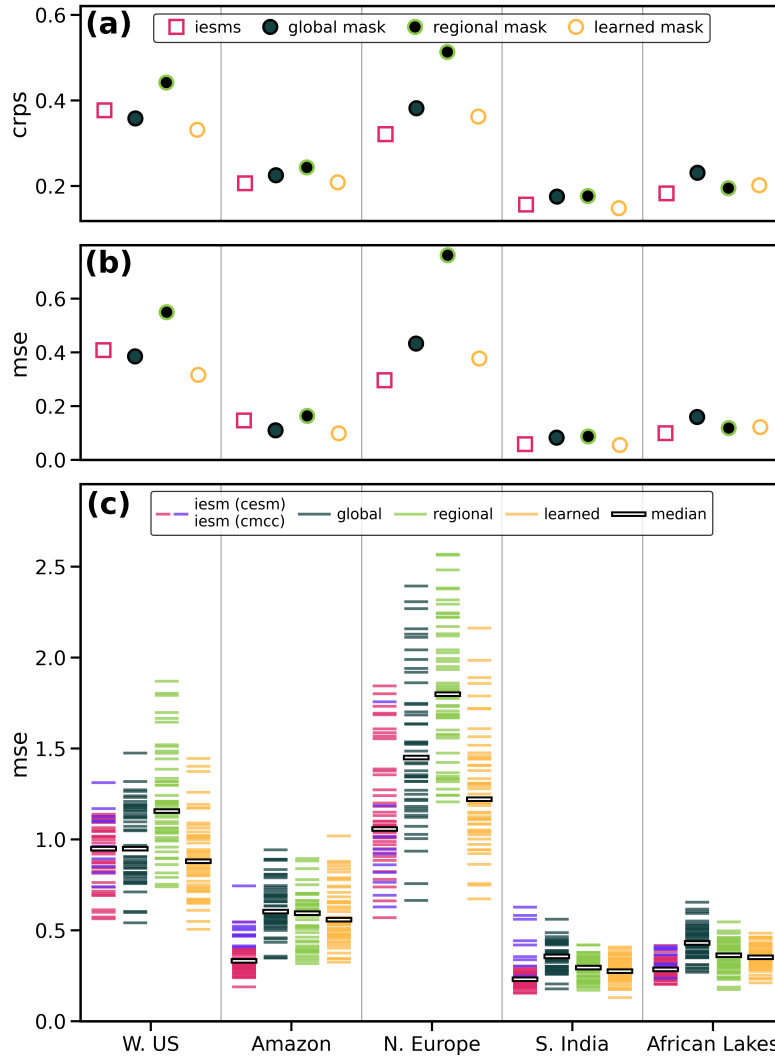


**Figure 10.** Same as Figure 9, but for the period 1999-2018. Bias-corrected IESMs (CESM and CMCC, see Table 2) can be directly compared for this time period.

Figure 10 includes the IESMs, and covers the twenty year time period 1999-2008. Once again, the learned mask is as good or better than the other analog methods for all regions and metrics. The learned mask outperforms the IESMs in the western United States region. The learned mask performs similarly to the IESMs in the Amazon, southern India, and the African Great Lakes, while the IESMs outperform the learned mask in northern Europe. For the analog forecast MSE, the learned mask outperforms the IESMs in the western United States, while the IESMs outperform the learned masks in the Amazon, northern Europe and the African Great Lakes.

## 4 Conclusions

We have presented a framework for multi-year-to-decadal prediction of regional 2-meter temperature. The framework incorporates a machine learned mask of weights that highlights important precursors for use in model-analog forecasts. Using several metrics, we have shown that our method:

- better captures the mean regional trend (ensemble mean analog MSE) and ensemble distribution (CRPS) at all lead times and regions compared to alternative analog methods, and lead times greater than two years and three of five regions compared to two initialized Earth system models (IESMs).
- better captures regional variability (analog forecast MSE) at all lead times and four of five regions (not the Amazon) compared to alternative analog methods, and lead times greater than two years compared to IESMs.
- better captures temperature patterns (EMD, Supplementary Material) at all lead times and regions compared to alternative analog methods, and about half the lead times and regions compared to IESMs.
- produces better classification (4 classes, Supplementary Material) for all lead times and regions compared to alternative analog methods, and lead times greater than two years and three of five regions compared to IESMs.

Our machine learning model-analog framework can be further improved. In this work, we only explored matching on and predicting the same variable, though multiple matching variables is likely to be beneficial. We used tethering (matching on multiple years) and explored a learned mask that was refined using transfer learning, but did not optimize either of these processes. In addition to improvements, the analog framework produces ensembles that can be used for prediction of extremes in a probabilistic sense.

There are many benefits to multi-year-to-decadal model-analog predictions. Analog predictions do not require expensive model runs, beyond what is currently available, and benefit from improved modeling (both in quality of analogs and number available). Analog predictions avoid the issue of initialization shock and climate drift faced by initialzaed Earth system models, reducing the need for bias corrections (Meehl et al., 2021, 2022), though our method still assumes that the models evolve consistently with observations. Furthermore, the mask of precursor weights makes this prediction framework interpretable. For example, Rader and Barnes (2023) calculate how much skill is attributable to known precursors by isolating/occluding those precursors. Analyses such as this are a strength of our interpretable approach, allowing us to identify predictable signals/patterns and connect them to physical understanding.

## Open Research Section

Berkeley Earth Surface Temperature (BEST; Rohde & Hausfather, 2020) can be obtained from `https://berkeleyearth.org/data/`, CMIP6 data and initialized Earth system model runs can be obtained from the Earth System Federation Grid (Cinquini et al., 2014) at `https://esgf.github.io/index.html`. Code underlying the machine learning model-analog framework presented will be made available upon publication at `https://github.com/mafern/mask-analog-predictions`.

# References

Acosta Navarro, J. C., Aranyossy, A., De Luca, P., Donat, M. G., Hrast Essen-felder, A., Mahmood, R., . . . Volpi, D. (2025). Seamless seasonal to multi-annual predictions of temperature and standardized precipitation index by constraining transient climate model simulations. *EGUsphere*, *2025*, 1–24. Retrieved from `https://egusphere.copernicus.org/preprints/2025/egusphere-2025-319/` doi: 10.5194/egusphere-2025-319

Befort, D. J., Brunner, L., Borchert, L. F., O'Reilly, C. H., Mignot, J., Ballinger, A. P., . . . Weisheimer, A. (2022). Combination of decadal predictions and climate projections in time: Challenges and potential solutions. *Geophysical Research Letters*, *49*(15), e2022GL098568. Retrieved from `https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2022GL098568` (e2022GL098568 2022GL098568) doi: https://doi.org/10.1029/2022GL098568

Befort, D. J., O'Reilly, C. H., & Weisheimer, A. (2020). Constraining projections using decadal predictions. *Geophysical Research Letters*, *47*(18), e2020GL087900. Retrieved from `https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2020GL087900` (e2020GL087900 10.1029/2020GL087900) doi: https://doi.org/10.1029/2020GL087900

Bergen, R. E., & Harnack, R. P. (1982). Long-range temperature prediction using a simple analog approach. *Monthly Weather Review*, *110*(8), 1083 - 1099. Retrieved from `https://journals.ametsoc.org/view/journals/mwre/110/8/1520-0493_1982_110_1083_lrtpua_2_0_co_2.xml` doi: 10.1175/1520-0493(1982)110⟨1083:LRTPUA⟩2.0.CO;2

Cinquini, L., Crichton, D., Mattmann, C., Harney, J., Shipman, G., Wang, F., . . . Schweitzer, R. (2014). The earth system grid federation: An open infrastructure for access to distributed geospatial data. *Future Generation Computer Systems*, *36*, 400-417. Retrieved from `https://www.sciencedirect.com/science/article/pii/S0167739X13001477` (Special Section: Intelligent Big Data Processing Special Section: Behavior Data Security Issues in Network Information Propagation Special Section: Energy-efficiency in Large Distributed Computing Architectures Special Section: eScience Infrastructure and Applications) doi: https://doi.org/10.1016/j.future.2013.07.002

De Luca, P., Delgado-Torres, C., Mahmood, R., Samso-Cabre, M., & Donat, M. G. (2023, sep). Constraining decadal variability regionally improves near-term projections of hot, cold and dry extremes. *Environmental Research Letters*, *18*(9), 094054. Retrieved from `https://dx.doi.org/10.1088/1748-9326/acf389` doi: 10.1088/1748-9326/acf389

Ding, H., & Alexander, M. A. (2023). Multi-year predictability of global sea surface temperature using model-analogs. *Geophysical Research Letters*, *50*(21), e2023GL104097. Retrieved from `https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2023GL104097` (e2023GL104097 2023GL104097) doi: https://doi.org/10.1029/2023GL104097

Ding, H., Newman, M., Alexander, M. A., & Wittenberg, A. T. (2018). Skillful climate forecasts of the tropical indo-pacific ocean using model-analogs. *Journal of Climate*, *31*(14), 5437 - 5459. Retrieved from `https://journals.ametsoc.org/view/journals/clim/31/14/jcli-d-17-0661.1.xml` doi: 10.1175/JCLI-D-17-0661.1

Donat, M. G., Mahmood, R., Cos, P., Ortega, P., & Doblas-Reyes, F. (2024, jun). Improving the forecast quality of near-term climate projections by constraining internal variability based on decadal predictions and observations. *Environmental Research: Climate*, *3*(3), 035013. Retrieved from `https://dx.doi.org/10.1088/2752-5295/ad5463` doi: 10.1088/2752-5295/ad5463

Dunstone, N., Lockwood, J., Solaraju-Murali, B., Reinhardt, K., Tsartsali, E. E., Athanasiadis, P. J., . . . Thornton, H. E. (2022). Towards useful decadal climate services. *Bulletin of the American Meteorological Society*, *103*(7), E1705

- E1719. Retrieved from `https://journals.ametsoc.org/view/journals/bams/103/7/BAMS-D-21-0190.1.xml` doi: 10.1175/BAMS-D-21-0190.1

Eyring, V., Bony, S., Meehl, G. A., Senior, C. A., Stevens, B., Stouffer, R. J., & Taylor, K. E. (2016). Overview of the coupled model intercomparison project phase 6 (cmip6) experimental design and organization. *Geoscientific Model Development*, *9*(5), 1937–1958. Retrieved from `https://gmd.copernicus.org/articles/9/1937/2016/` doi: 10.5194/gmd-9-1937-2016

Ghani, B., Kalkman, V. J., Planqué, B., Vellinga, W.-P., Gill, L., & Stowell, D. (2024). Generalization in birdsong classification: impact of transfer learning methods and dataset characteristics. *arXiv e-prints*, arXiv:2409.15383. doi: 10.48550/arXiv.2409.15383

Gneiting, T., & Raftery, A. E. (2007). Strictly proper scoring rules, prediction, and estimation. *Journal of the American Statistical Association*, *102*(477), 359–378. Retrieved from `https://doi.org/10.1198/016214506000001437` doi: 10.1198/016214506000001437

Hermanson, L., Smith, D., Seabrook, M., Bilbao, R., Doblas-Reyes, F., Tourigny, E., ... Kumar, A. (2022). Wmo global annual to decadal climate update: A prediction for 2021–25. *Bulletin of the American Meteorological Society*, *103*(4), E1117 - E1129. Retrieved from `https://journals.ametsoc.org/view/journals/bams/103/4/BAMS-D-20-0311.1.xml` doi: 10.1175/BAMS-D-20-0311.1

IPCC. (2023a). *Climate change 2021 – the physical science basis: Working group i contribution to the sixth assessment report of the intergovernmental panel on climate change.* Cambridge University Press.

IPCC. (2023b). *Climate change 2022 - mitigation of climate change: Working group iii contribution to the sixth assessment report of the intergovernmental panel on climate change.* Cambridge University Press.

IPCC. (2023c). *Climate change 2022 – impacts, adaptation and vulnerability: Working group ii contribution to the sixth assessment report of the intergovernmental panel on climate change.* Cambridge University Press.

Khasnis, A. A., & Nettleman, M. D. (2005). Global warming and infectious disease. *Archives of Medical Research*, *36*(6), 689-696. Retrieved from `https://www.sciencedirect.com/science/article/pii/S0188440905001517` (Infectious Diseases: Revisiting Past Problems and Addressing Future Challenges) doi: https://doi.org/10.1016/j.arcmed.2005.03.041

Kim, Y.-Y., Lee, J.-Y., Timmermann, A., Chikamoto, Y., Lee, S.-S., Kwon, E. Y., ... Franzke, C. (2025). Robust estimates of earth system predictability of the 1st kind using the cesm2 multiyear prediction system (cesm2-mp). *Research Square*. Retrieved from `https://www.researchsquare.com/article/rs-5748726/v1` doi: 10.21203/rs.3.rs-5748726/v1

Lorenz, E. N. (1969, January). Atmospheric Predictability as Revealed by Naturally Occurring Analogues. *Journal of the Atmospheric Sciences*, *26*(4), 636. doi: 10.1175/1520-0469(1969)26⟨636:APARBN⟩2.0.CO;2

Mahmood, R., Donat, M. G., Ortega, P., Doblas-Reyes, F. J., & Ruprich-Robert, Y. (2021). Constraining decadal variability yields skillful projections of near-term climate change. *Geophysical Research Letters*, *48*(24), e2021GL094915. Retrieved from `https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2021GL094915` (e2021GL094915 2021GL094915) doi: https://doi.org/10.1029/2021GL094915

Meehl, G. A., Richter, J. H., Teng, H., Capotondi, A., Cobb, K., Doblas-Reyes, F., ... Xie, S.-P. (2021). Initialized earth system prediction from subseasonal to decadal timescales. *Nature Reviews Earth and Environment*, *2*(5), 340-357. doi: 10.1038/s43017-021-00155-x

Meehl, G. A., Teng, H., Smith, D., Yeager, S., Merryfield, W., Doblas-Reyes, F., & Glanville, A. A. (2022). The effects of bias, drift, and trends in

calculating anomalies for evaluating skill of seasonal-to-decadal initialized climate predictions. *Climate Dynamics*, *59*(11-12), 3373-3389. doi: 10.1007/s00382-022-06272-7

Pan, S. J., & Yang, Q. (2010). A survey on transfer learning. *IEEE Transactions on Knowledge and Data Engineering*, *22*(10), 1345-1359. doi: 10.1109/TKDE .2009.191

Pele, O., & Werman, M. (2008, October). A linear time histogram metric for improved sift matching. In *Computer vision–eccv 2008* (pp. 495–508). Springer.

Pele, O., & Werman, M. (2009, September). Fast and robust earth mover's distances. In *2009 ieee 12th international conference on computer vision* (pp. 460–467).

Rader, J. K., & Barnes, E. A. (2023). Optimizing seasonal-to-decadal analog forecasts with a learned spatially-weighted mask. *Geophysical Research Letters*, *50*(23), e2023GL104983. Retrieved from `https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2023GL104983` (e2023GL104983 2023GL104983) doi: https://doi.org/10.1029/2023GL104983

Rohde, R. A., & Hausfather, Z. (2020). The berkeley earth land/ocean temperature record. *Earth System Science Data*, *12*(4), 3469–3479. Retrieved from `https://essd.copernicus.org/articles/12/3469/2020/` doi: 10.5194/essd -12-3469-2020

Rubner, Y., Tomasi, C., & Guibas, L. (1998). A metric for distributions with applications to image databases. In *Sixth international conference on computer vision (ieee cat. no.98ch36271)* (p. 59-66). doi: 10.1109/ICCV.1998.710701

Solaraju-Murali, B., Bojovic, D., Gonzalez-Reviriego, N., Nicodemou, A., Terrado, M., Caron, L.-P., & Doblas-Reyes, F. J. (2022). How decadal predictions entered the climate services arena: an example from the agriculture sector. *Climate Services*, *27*, 100303. Retrieved from `https://www.sciencedirect.com/science/article/pii/S2405880722000218` doi: https://doi.org/10.1016/j.cliser.2022.100303

Toride, K., Newman, M., Hoell, A., Capotondi, A., Schlör, J., & Amaya, D. J. (2024). Using Deep Learning to Identify Initial Error Sensitivity for Interpretable ENSO Forecasts. *arXiv e-prints*, arXiv:2404.15419. doi: 10.48550/arXiv.2404.15419

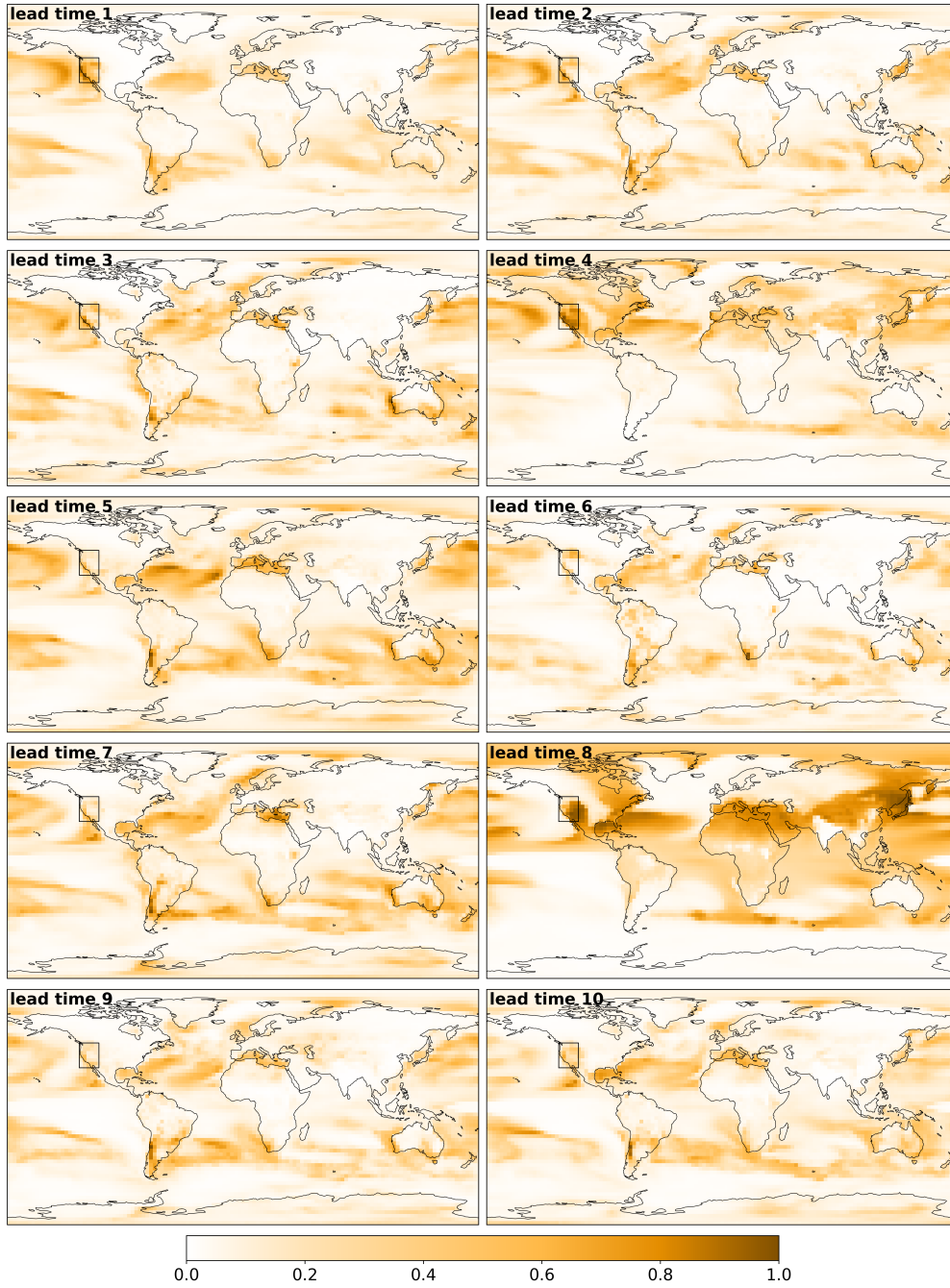# Appendix A  Supplementary Material



**Figure A1.**  Model learned masks for the western United States for lead times 1-10 years.
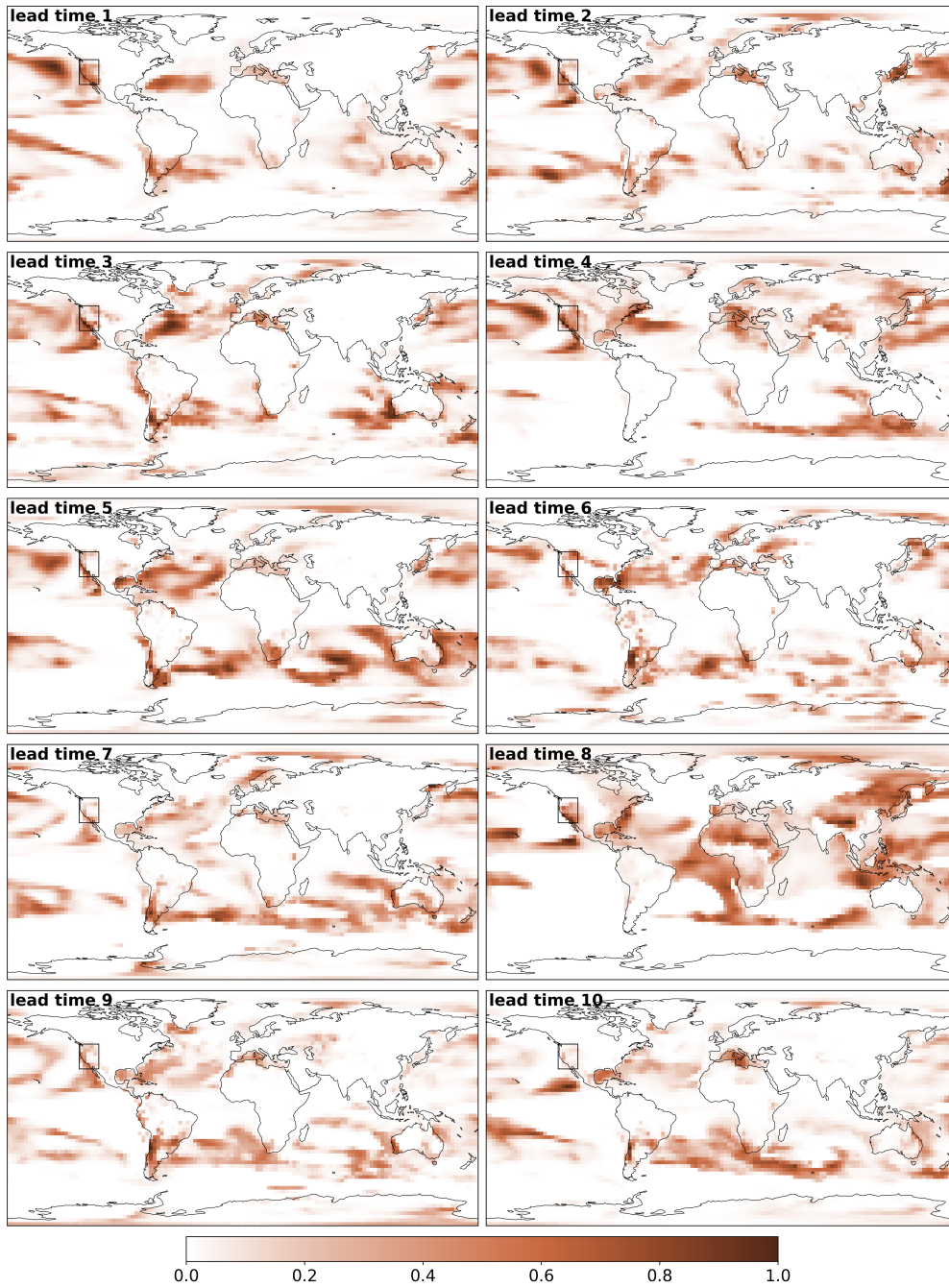
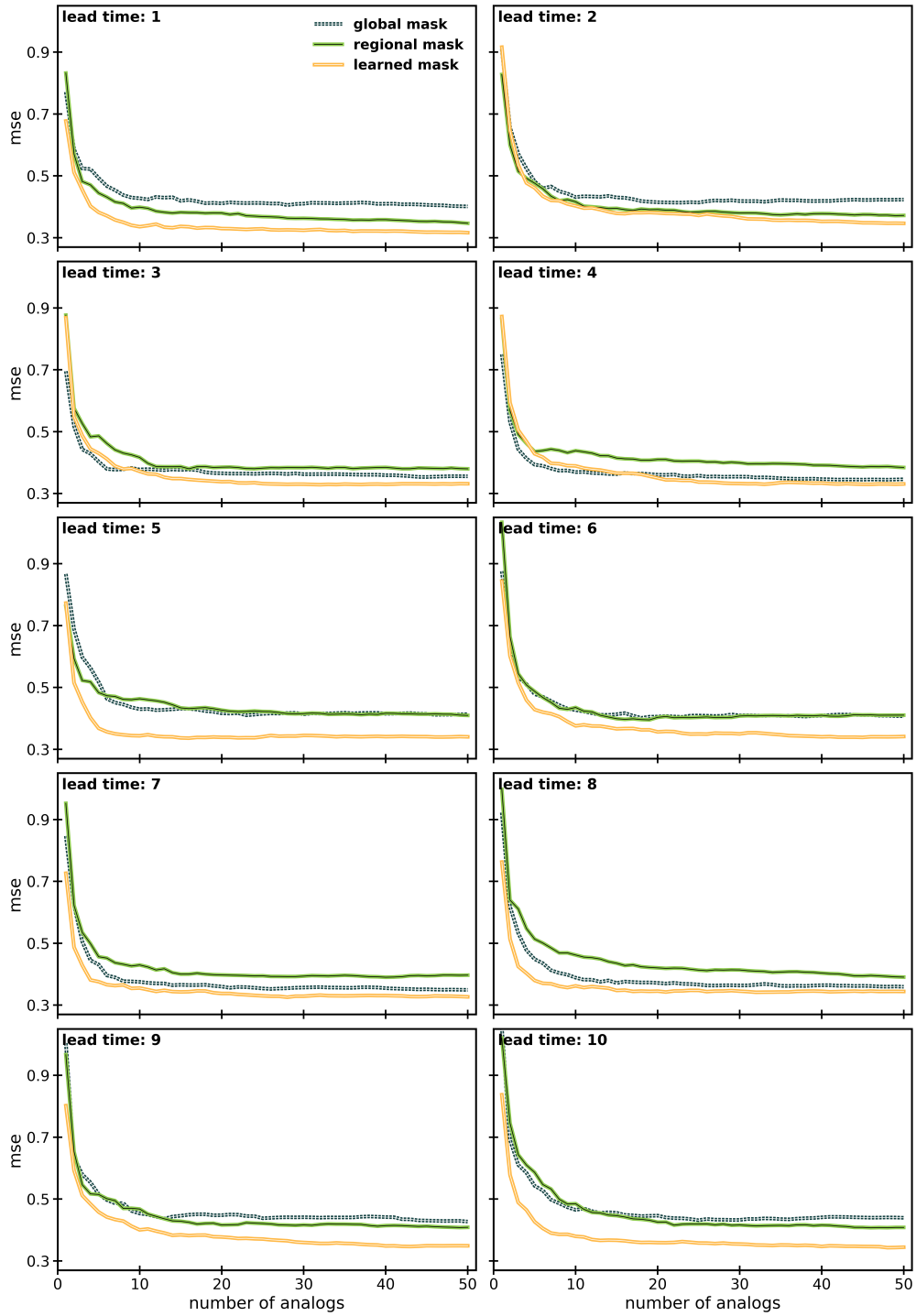**Figure A2.** Transfer learned masks for the western United States for lead times 1-10 years.

**Figure A3.** MSE versus number analogs used to calculate the mean prediction, for all lead times. Also shown are the global and regional mask results.
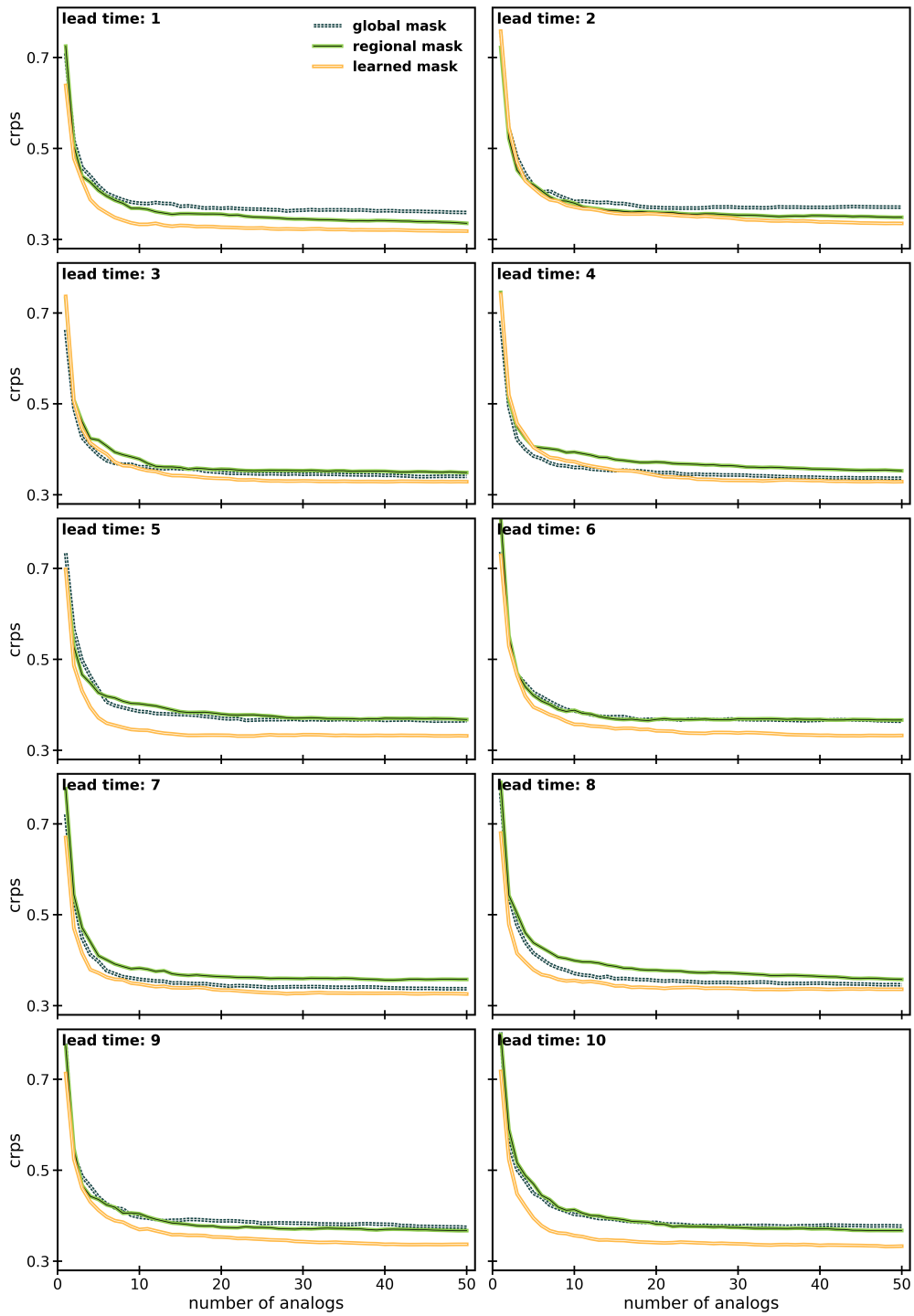
**Figure A4.** Same as Figure A3, but for CRPS.

**Figure A5.** EMD versus lead time for the western United States for mean predictions (line plots) and individual analogs (distributions). The top two panels show EMD for the time period 1864-2023, while the bottom two panels show EMD for 2009-2018.
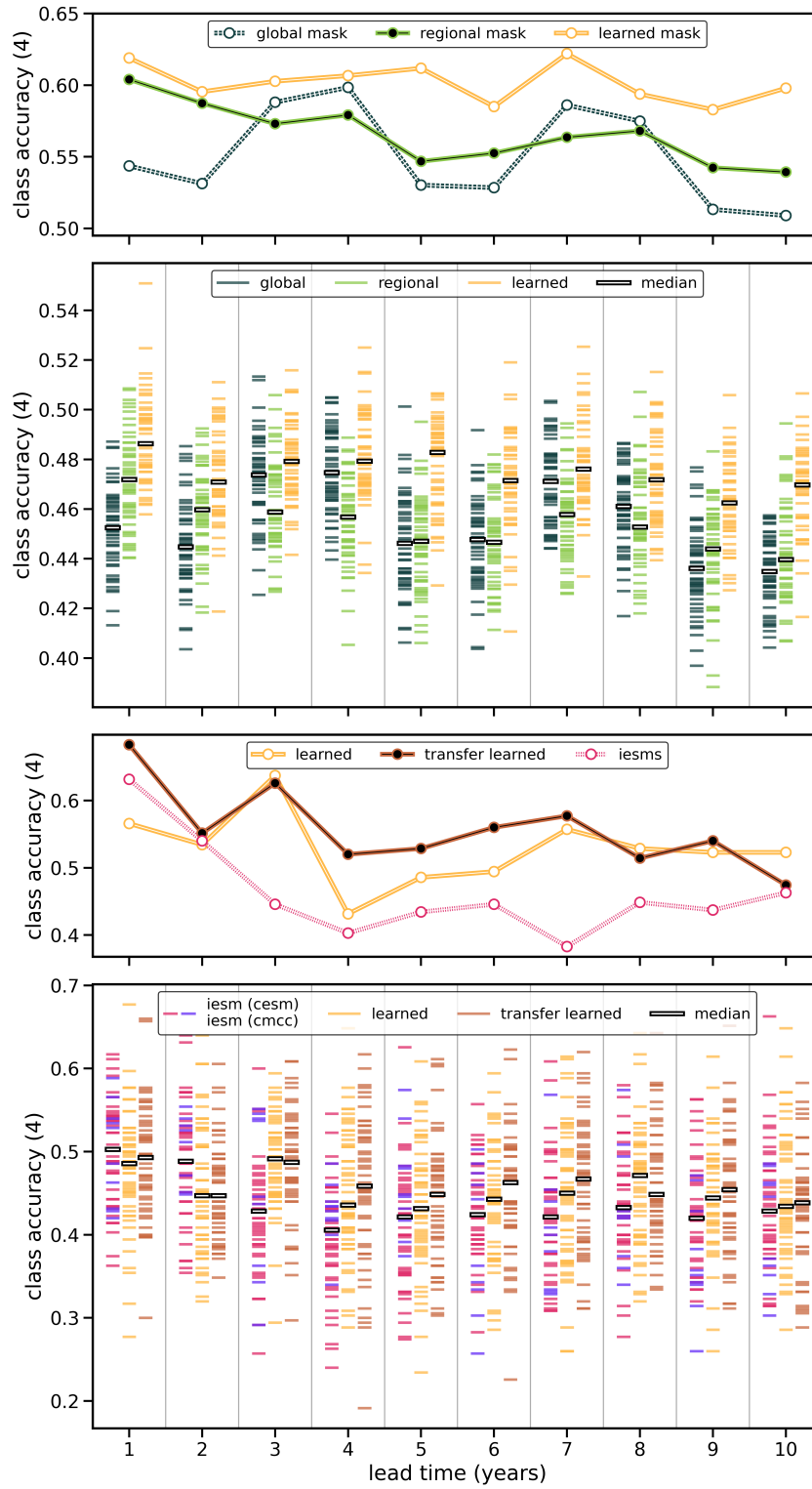
**Figure A6.** Same as Figure A5, but for class accuracy, with four classes.
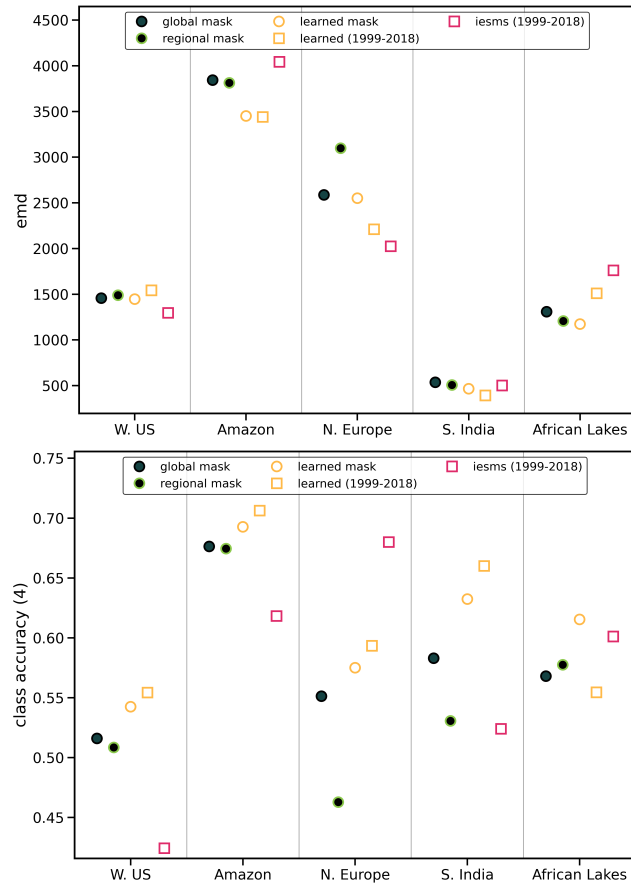
**Figure A7.** EMD (top) and class accuracy (bottom, 4 classes). The circles are the mean metric covering the time period 1956-2023, the squares cover 1999-2018.
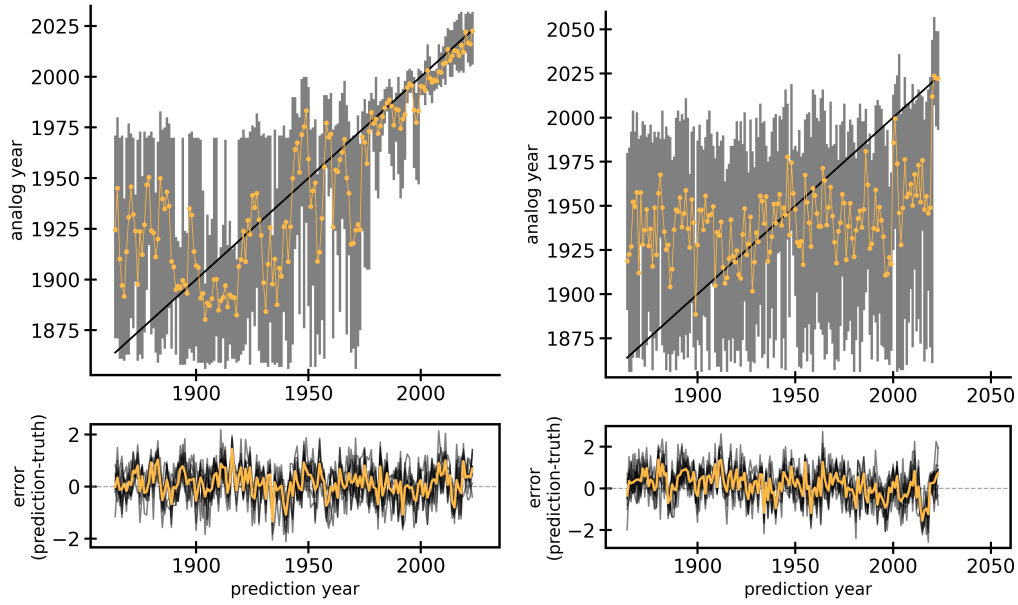
**Figure A8.** For each prediction year, the range of years for the top ten analogs (top panel) as well as the error over that time (bottom). While the model learned mask (left) begins following the one-to-one line (indicating that the selected analog year is the same as the prediction year) by the 1970s, the regional mask (right) only does so in the final few years. In the bottom panels, it can be seen that the model learned mask predictions are not biased toward underestimating or overestimating at any point, while the regional mask predictions tends towards underestimation starting in the 1970s. This is likely due to the regional mask not being able to pick up on the forced warming signal, which may be most apparent outside the target region.