

716 A Licenses and Terms of Use

717 ClimateLearn is a software package that can be installed from the Python Package Index as follows.

```
pip install climate-learn
```

718 The source code is available online under the MIT License at [https://github.com/](https://github.com/aditya-grover/climate-learn)
719 [aditya-grover/climate-learn](https://github.com/aditya-grover/climate-learn), and the accompanying documentation website is at <https://climatelearn.readthedocs.io/>. The Extreme-ERA5 dataset does not exist as a distinct entity, but can be produced by running code provided in our library. The Machine Intelligence Group at
721 UCLA is the maintainer of ClimateLearn.
722

723 The sources for datasets provided by ClimateLearn are WeatherBench, ClimateBench, the Earth
724 System Grid Federation (ESGF), the Copernicus Climate Data Store (CDS), and PRISM. The
725 WeatherBench dataset (<https://mediatum.ub.tum.de/1524895>), ClimateBench dataset (<https://zenodo.org/record/7064308>), and MPI-ESM1.2-HR outputs from ESGF (<https://pcmdi.llnl.gov/CMIP6/TermsOfUse/TermsOfUse6-2.html>) are available under the CC BY 4.0 license. Neither Copernicus (<https://cds.climate.copernicus.eu/api/v2/terms/static/licence-to-use-copernicus-products.pdf>) nor PRISM (<https://prism.oregonstate.edu/terms/>) use a Creative Commons License. Instead, they each set forth their own terms
729 of use, both of which permit reproduction and distribution for non-commercial purposes.
730
731

732 B Experiment details

733 B.1 Network architectures

734 B.1.1 ResNet

735 Our ResNet architecture is similar to that of WeatherBench [61, 60], in which each residual block
736 consists of two identical convolutional modules: 2D convolution \rightarrow LeakyReLU with $\alpha = 0.3 \rightarrow$
737 Batch Normalization \rightarrow Dropout.

Table 4: Default hyperparameters of ResNet

Hyperparameter	Meaning	Value
Padding size	Padding size of each convolution layer	1
Kernel size	Kernel size of each convolution layer	3
Stride	Stride of each convolution layer	1
Hidden dimension	Number of output channels of each residual block	128
Residual blocks	Number of residual blocks	28
Dropout	Dropout rate	0.1

738 Table 4 shows the hyperparameters for ResNet in all of our experiments. We use a convolutional layer
739 with a kernel size of 7 at the beginning of the network. All paddings are periodic in the longitude
740 direction and zeros in the latitude direction.

741 B.1.2 UNet

742 We borrow our UNet implementation from PDEArena [21]. Table 5 shows the hyperparameters for
743 UNet in all of our experiments. Similar to ResNet, we use a convolutional layer with a kernel size of
744 7 at the beginning of the network, and all paddings are periodic in the longitude direction and zeros
745 in the latitude direction.

Table 5: Default hyperparameters of UNet

Hyperparameter	Meaning	Value
Padding size	Padding size of each convolution layer	1
Kernel size	Kernel size of each convolution layer	3
Stride	Stride of each convolution layer	1
Hidden dimension	Base number of output channels	64
Channel multiplications	Determine the number of output channels for Down and Up blocks	[1, 2, 2]
Blocks	Number of blocks	2
Use attention	If use attention in Down and Up blocks	False
Dropout	Dropout rate	0.1

746 B.1.3 ViT

747 We use the standard Vision Transformer architecture [14] with minor modifications. We remove
 748 the class token and add a 1-hidden MLP prediction head which is applied to the tokens after the
 749 last attention layer to predict the outputs. Tabel 6 shows the hyperparameters for ViT in all of our
 experiments.

Table 6: Default hyperparameters of ViT

Hyperparameter	Meaning	Value
p	Patch size	2
D	Embedding dimension	128
Depth	Number of ViT blocks	8
# heads	Number of attention heads	4
MLP ratio	Determine the hidden dimension of the MLP layer in a ViT block	4
Prediction depth	Number of layers of the prediction head	2
Hidden dimension	Hidden dimension of the prediction head	128
Drop path	For stochastic depth [30]	0.1
Dropout	Dropout rate	0.1

750

751 B.2 Datasets

752 B.2.1 ERA5

753 We refer to <https://confluence.ecmwf.int/display/CKB/ERA5%3A+data+documentation>
 754 for more details of the raw ERA5 data. We use the preprocessed version of ERA5 at 5.625° from
 755 WeatherBench [61] for our experiments. Table 7 summarizes the variables we use for our experiments.
 756

757 B.2.2 Extreme-ERA5

758 **Calculating thresholds** We use the surface temperature (T2m) data corresponding to the years
 759 1979 – 2015 from ERA5 at a resolution of 5.625° to calculate the thresholds. The thresholds are
 760 localized i.e. they are calculated for every pixel on the grid. For a given timestamp and pixel, we
 761 first calculate a 7 day mean till that timestamp. Now, to account for neighboring regions/pixels, we
 762 set the localized mean as $0.44 * \text{current pixel's mean} + 0.11 * \text{sum of means of pixels sharing an}$
 763 $\text{edge} + 0.027 * \text{sum of means of pixels sharing a vertex but not an edge}$. Note, that there is no need
 764 of padding while accounting for neighboring pixels, since earth is a globe, the neighbors of leftmost
 765 pixels include the rightmost pixels and vice-versa. Finally, the 5th and 95th percentile values of this
 766 new mean data corresponding to every pixel is set as threshold.

Table 7: ERA5 variables used in our experiments. *Constant* represents constant variables, *Single* represents surface variables, and *Atmospheric* represents atmospheric properties at the chosen altitudes.

Type	Variable name	Abbrev.	Levels
Static	Land-sea mask	LSM	
Static	Orography		
Static	Latitude		
Single	Toa incident solar radiation	Tsr	
Single	2 metre temperature	T2m	
Single	10 metre U wind component	U10	
Single	10 metre V wind component	V10	
Atmospheric	Geopotential	Z	50, 250, 500, 600, 700, 850, 925
Atmospheric	U wind component	U	50, 250, 500, 600, 700, 850, 925
Atmospheric	V wind component	V	50, 250, 500, 600, 700, 850, 925
Atmospheric	Temperature	T	50, 250, 500, 600, 700, 850, 925
Atmospheric	Specific humidity	Q	50, 250, 500, 600, 700, 850, 925
Atmospheric	Relative humidity	R	50, 250, 500, 600, 700, 850, 925

767 **Building masks** As the purpose of Extreme-ERA5 is evaluation of forecasting models under extreme
768 weather conditions, we build it for test years i.e. 2017 – 2018 only. We first create a 2-D mask of
769 size, latitude x longitude, filled with zeros for every available timestamp in the test years. Similar to
770 the calculating thresholds, we compute the mean of each pixel at every timestamp for T2m’s test data.
771 We then, set the value for a given pixel in the mask as 1, if the mean value is outside the bounds set
772 by the thresholds. Finally, during evaluation time, we use these masks to select subset of data.

773 B.2.3 CMIP6

774 **MPI-ESM1.2-HR** We use MPI-ESM1.2-HR, a dataset in the CMIP6 data repository for our
experiments in Section 4.1.3. Table 8 summarizes the variables we use for our experiments.

Table 8: MPI-ESM1.2-HR variables used in our experiments. *Single* represents surface variables and *Atmospheric* represents atmospheric properties at the chosen altitudes.

Type	Variable name	Abbrev.	Levels
Single	2 metre temperature	T2m	
Single	10 metre U wind component	U10	
Single	10 metre V wind component	V10	
Atmospheric	Geopotential	Z	50, 250, 500, 600, 700, 850, 925
Atmospheric	U wind component	U	50, 250, 500, 600, 700, 850, 925
Atmospheric	V wind component	V	50, 250, 500, 600, 700, 850, 925
Atmospheric	Temperature	T	50, 250, 500, 600, 700, 850, 925
Atmospheric	Specific humidity	Q	50, 250, 500, 600, 700, 850, 925

775

776 **ClimateBench** We adopt data from ClimateBench [82] for our climate projection experiment.
777 ClimateBench contains simulated data from experimental runs by the Norwegian Earth System
778 Model [69], a member of CMIP6, on different emission scenarios. Specifically, ClimateBench
779 includes 7 emission scenarios: historical, ssp126, ssp370, ssp585, hist-aer, hist-GHG, and ssp245.
780 We refer to the original ClimateBench paper for the exact temporal coverage and more details of
781 these scenarios.

782 **B.3 Training details**

783 **B.3.1 Continuous training**

784 Continuous models additionally condition on lead times to make predictions. To do this, we add
 785 the lead time value in hours divided by 100 to the input channels to make the model aware of the
 786 lead time it is forecasting at. During training, we randomize the lead time from 6 hours to 5 days
 787 $\Delta t \sim \mathcal{U}[6, 120]$, and during evaluation, we fix the lead time to a certain value to evaluate the model’s
 788 performance at a certain lead time. This setting was commonly used in previous works [50, 61].

789 **B.3.2 Software and hardware stack**

790 We use PyTorch [51], numpy [24] and xarray [29] to manage our data and model training. We also
 791 use timm [86] for our ViT implementation. All training is done on 10 AMD EPYC 7313 CPU cores
 792 and one NVIDIA RTX A5000 GPU. We leverage fp16 floating point precision in our experiments.

793 **B.4 Metrics**

794 We use the following definitions in our metric formulations

- 795 • N is the number of data points
- 796 • H is the number of latitude coordinates.
- 797 • W is the number of longitude coordinates.
- 798 • X and \tilde{X} are the ground-truth and prediction, respectively.

799 The latitude weighting function is given by

$$L(i) = \frac{\cos(H_i)}{\frac{1}{H} \sum_{i=1}^H \cos(H_i)} \quad (1)$$

800 **B.4.1 Deterministic weather forecasting metrics**

801 **Root mean square error (RMSE)**

$$\text{RMSE} = \frac{1}{N} \sum_{k=1}^N \sqrt{\frac{1}{H \times W} \sum_{i=1}^H \sum_{j=1}^W L(i) (\tilde{X}_{k,i,j} - X_{k,i,j})^2}. \quad (2)$$

802 **Anomaly correlation coefficient (ACC)** is the spatial correlation between prediction anomalies \tilde{X}'
 803 relative to climatology and ground truth anomalies X' relative to climatology:

$$\text{ACC} = \frac{\sum_{k,i,j} L(i) \tilde{X}'_{k,i,j} X'_{k,i,j}}{\sqrt{\sum_{k,i,j} L(i) \tilde{X}'_{k,i,j}^2 \sum_{k,i,j} L(i) X'_{k,i,j}^2}}, \quad (3)$$

$$\tilde{X}' = \tilde{X} - C, X' = X - C, \quad (4)$$

804 in which climatology C is the temporal mean of the ground truth data over the entire test set
 805 $C = \frac{1}{N} \sum_k X$.

806 **B.4.2 Probabilistic weather forecasting metrics**

807 **Spread-skill ratio (Spread by RMSE)** measures a probabilistic forecast’s reliability. Let N be the
 808 number of forecasts produced either by ensembling or drawing samples from a parametric prediction.
 809 Spread is given by

$$\text{Spread} = \frac{1}{N} \sum_k \sqrt{\frac{1}{H \times W} \sum_{i=1}^H \sum_{j=1}^W L(i) \text{var}(\tilde{X}_{i,j})} \quad (5)$$

810 **Continuous ranked probability score** measures a probabilistic forecast's calibration and sharpness.
 811 Let F denote the CDF of the forecast distribution. For a Gaussian distribution parameterized by mean
 812 μ and standard deviation σ , the closed-form, differentiable solution is

$$\text{CRPS}(F_{\mu,\sigma}, X) = \sigma \left\{ \frac{X - \mu}{\sigma} \left[2\Phi \left(\frac{X - \mu}{\sigma} \right) - 1 \right] + 2\phi \left(\frac{X - \mu}{\sigma} \right) - \frac{1}{\sqrt{\pi}} \right\} \quad (6)$$

813 where Φ and ϕ are the CDF and PDF of the standard normal distribution, respectively.

814 B.4.3 Climate downscaling metrics

815 **Root mean square error (RMSE)** This is the same as Equation (2).

816 **Mean bias** measures the difference between the spatial mean of the prediction and the spatial mean
 817 of the ground truth. A positive mean bias shows an overestimation, while a negative mean bias shows
 818 an underestimation of the mean value.

$$\text{Mean bias} = \frac{1}{N \times H \times W} \sum_{k=1}^N \sum_{i=1}^H \sum_{j=1}^W \tilde{X} - \frac{1}{N \times H \times W} \sum_{k=1}^N \sum_{i=1}^H \sum_{j=1}^W X \quad (7)$$

819 **Pearson coefficient** measures the correlation between the prediction and the ground truth. We first
 820 flatten the prediction and ground truth, and compute the metric as follows:

$$\rho_{\tilde{X}, X} = \frac{\text{cov}(\tilde{X}, X)}{\sigma_{\tilde{X}} \sigma_X} \quad (8)$$

821 **Masking for PRISM** Since PRISM does not record data over the oceans, we mask out those values
 822 for evaluation. Concretely, we set NaN values in the ground truth data to 0. Then, we multiply the
 823 model's predictions by a binary mask that is 0 wherever the ground truth data is originally NaN and is
 824 1 everywhere else.

825 B.4.4 Climate projection metrics

826 **Normalized spatial root mean square error (NRMSE_s)** measures the spatial discrepancy between
 827 the temporal mean of the prediction and the temporal mean of the ground truth:

$$\text{NRMSE}_s = \sqrt{\left\langle \left(\frac{1}{N} \sum_{k=1}^N \tilde{X} - \frac{1}{N} \sum_{k=1}^N X \right)^2 \right\rangle / \frac{1}{N} \sum_{k=1}^N \langle X \rangle}, \quad (9)$$

828 in which $\langle A \rangle$ is the global mean of A :

$$\langle A \rangle = \frac{1}{H \times W} \sum_{i=1}^H \sum_{j=1}^W L(i) A_{i,j} \quad (10)$$

829 **Normalized global root mean square error (NRMSE_g)** measures the discrepancy between the
 830 global mean of the prediction and the global mean of the ground truth:

$$\text{NRMSE}_g = \sqrt{\frac{1}{N} \sum_{k=1}^N \left(\langle \tilde{X} \rangle - \langle X \rangle \right)^2 / \frac{1}{N} \sum_{k=1}^N \langle X \rangle}. \quad (11)$$

831 **Total normalized root mean square error (Total)** is the weighted sum of NRMSE_s and NRMSE_g:

$$\text{Total} = \text{NRMSE}_s + \alpha \cdot \text{NRMSE}_g, \quad (12)$$

832 where α is chosen to be 5 as suggested by Watson-Parris et al. [82].

833 **C Additional experiments**

834 **C.1 Climate projection**

835 **Task** We consider the task of predicting the annual mean distributions of 4 target variables in
 836 ClimateBench [82]: surface temperature, diurnal temperature range, precipitation, and the 90th
 837 percentile of precipitation.

838 **Baselines** We compare ResNet, UNet, and ViT, three deep learning models supported by
 839 ClimateLearn with CNN-LSTM, the deep learning baseline in ClimateBench. The network archi-
 840 tectures of the three models are identical to Appendix B.1.

841 **Data** We regrid the original ClimateBench data to 5.625° for easy training and evaluation. The input
 842 variables include 4 forcing factors: carbon dioxide (CO₂), sulfur dioxide (SO₂), black carbon (BC),
 843 and methane (CH₄). Similar to the deep learning baseline in ClimateBench, we stack 10 consecutive
 844 years to predict the target variables of the current year. We standardize the input channels to have 0
 845 mean and 1 standard deviation, but do not standardize the output variables. Training and validation
 846 data includes the historical data, ssp126, ssp370, ssp585, and the historical data with aerosol (hist-aer)
 847 and greenhouse gas (hist-GHG) forcings, and test data includes ssp245. We split train/validation data
 848 with a ratio of 0.9/0.1.

849 **Training and evaluation** We train one network for each target variable. We use the same optimizer
 850 and scheduler as in Section 4.1. We train for 50 epochs with 16 batch size, and use early stopping
 851 with a patience of 5 epochs. We use mean-squared error as the loss function and evaluation metric.
 852 We report normalized spatial root mean square error (NRMSE_s), normalized global root mean square
 853 error (NRMSE_g), and Total = NRMSE_s + 5 × NRMSE_g as test metrics.

854 **Results** Table 9 shows the performance of different baselines on ClimateBench. CNN-LSTM and
 855 UNet are the best-performing methods, with each achieving the best performance in 5/12 metrics,
 856 followed by ResNet which performs best on 2/12 metrics. ViT achieves a reasonable performance
 857 but underperforms the CNN-based methods.

Table 9: Performance of different deep learning baselines on ClimateBench. CNN-LSTM result is taken from ClimateBench.

	Surface temperature			Diurnal temperature range			Precipitation			90th percentile precipitation		
	NRMSE _s	NRMSE _g	Total	NRMSE _s	NRMSE _g	Total	NRMSE _s	NRMSE _g	Total	NRMSE _s	NRMSE _g	Total
CNN-LSTM	0.107	0.044	0.327	9.917	1.372	16.778	2.128	0.209	3.175	2.610	0.346	4.339
ResNet	0.182	0.042	0.395	9.128	0.737	12.810	2.930	0.180	3.828	3.413	0.286	4.845
UNet	0.097	0.046	0.328	6.300	0.946	11.030	2.483	0.141	3.187	3.122	0.282	4.532
ViT	0.191	0.092	0.650	7.725	0.746	11.460	2.909	0.327	4.545	3.615	0.418	5.704

858 **C.2 Extreme weather prediction**

859 Table 10 shows the performance of different models across various different lead times on the default
 860 test split and Extreme-ERA5. As discussed in Section 4.1.2, the performance of all models except
 861 Climatology is better on the extreme split than on the default split.

Table 10: Latitude-weighted RMSE on the normal and extreme test splits of ERA5 for different lead times.

T2M	6 Hours	1 Day	3 Days	5 Days	10 Days
Climatology	5.87 / 6.51	5.87 / 6.53	5.87 / 6.58	5.88 / 6.65	5.89 / 6.76
Persistence	2.76 / 2.99	2.13 / 1.78	2.99 / 2.42	3.26 / 2.61	3.59 / 2.89
ResNet	0.72 / 0.72	0.94 / 0.91	1.50 / 1.33	2.20 / 1.86	2.78 / 2.39
U-Net	0.76 / 0.77	1.04 / 0.99	1.65 / 1.43	2.26 / 1.88	2.76 / 2.44
ViT	0.78 / 0.80	1.09 / 1.05	1.71 / 1.55	2.38 / 2.04	2.78 / 2.30

Table 11: Performance of different models trained on one dataset (columns) and evaluated on another (rows). Training data for CMIP6 is available from the years 1850 – 2010, at a 6 hour frequency. Training data for ERA5 is available from years 1979 – 2010, at an one hour frequency.

		3 Days				5 Days				
		ERA5		CMIP6		ERA5		CMIP6		
		ACC	RMSE	ACC	RMSE	ACC	RMSE	ACC	RMSE	
ERA5	ResNet	Z500	0.95	315.07	0.96	302.08	0.77	646.57	0.86	531.47
		T850	0.93	1.84	0.91	2.08	0.80	3.00	0.83	2.77
		T2m	0.95	1.56	0.94	1.85	0.89	2.35	0.90	2.29
	U-Net	Z500	0.92	388.17	0.94	337.34	0.74	686.90	0.82	590.80
		T850	0.91	2.09	0.90	2.17	0.78	3.10	0.81	2.93
		T2m	0.95	1.72	0.93	1.89	0.89	2.38	0.89	2.37
	ViT	Z500	0.93	380.22	0.93	373.57	0.68	749.82	0.82	592.36
		T850	0.91	2.08	0.89	2.31	0.75	3.27	0.80	2.97
		T2m	0.94	1.73	0.92	2.10	0.88	2.54	0.88	2.52
CMIP6	ResNet	Z500	0.95	35.84	0.98	24.51	0.77	71.50	0.89	50.58
		T850	0.92	2.09	0.96	1.43	0.79	3.19	0.89	2.33
		T2m	0.94	1.88	0.97	1.32	0.88	2.54	0.94	1.87
	U-Net	Z500	0.92	43.36	0.96	30.61	0.75	74.68	0.85	58.67
		T850	0.90	2.30	0.95	1.67	0.78	3.29	0.87	2.57
		T2m	0.93	2.00	0.96	1.46	0.88	2.57	0.93	1.99
	ViT	Z500	0.93	42.19	0.95	34.83	0.68	83.68	0.85	58.86
		T850	0.90	2.25	0.94	1.83	0.75	3.48	0.86	2.60
		T2m	0.92	2.15	0.95	1.59	0.85	2.88	0.92	2.03

Table 12: Performance of different models trained on one dataset (columns) and evaluated on another (rows). The training years and data availability frequency is same for both the datasets.

		3 Days				5 Days				
		ERA5		CMIP6		ERA5		CMIP6		
		ACC	RMSE	ACC	RMSE	ACC	RMSE	ACC	RMSE	
ERA5	ResNet	Z500	0.95	322.86	0.94	345.00	0.79	624.20	0.78	646.48
		T850	0.93	1.90	0.90	2.21	0.81	2.91	0.79	3.11
		T2m	0.95	1.62	0.93	1.94	0.90	2.33	0.88	2.55
	U-Net	Z500	0.92	401.08	0.91	422.77	0.74	685.75	0.73	712.62
		T850	0.90	2.17	0.88	2.42	0.78	3.10	0.76	3.29
		T2m	0.94	1.81	0.91	2.19	0.89	2.44	0.86	2.73
	ViT	Z500	0.91	426.70	0.90	444.14	0.72	698.08	0.72	720.15
		T850	0.89	2.27	0.87	2.51	0.78	3.12	0.76	3.31
		T2m	0.94	1.88	0.91	2.19	0.89	2.43	0.87	2.69
CMIP6	ResNet	Z500	0.95	36.47	0.96	31.29	0.79	69.62	0.81	65.04
		T850	0.91	2.11	0.94	1.70	0.81	3.09	0.84	2.82
		T2m	0.93	1.91	0.96	1.53	0.88	2.51	0.91	2.24
	U-Net	Z500	0.92	44.95	0.93	40.98	0.74	75.45	0.76	72.67
		T850	0.89	2.37	0.92	2.05	0.78	3.27	0.81	3.04
		T2m	0.91	2.23	0.94	1.74	0.86	2.73	0.90	2.31
	ViT	Z500	0.91	46.92	0.92	43.91	0.73	76.80	0.75	74.22
		T850	0.89	2.40	0.91	2.15	0.77	3.29	0.80	3.06
		T2m	0.91	2.15	0.94	1.82	0.87	2.68	0.90	2.32

862 C.3 Dataset robustness

863 Table 11 shows the comparison of the performance for different models when trained on ERA5 and
864 evaluated on CMIP6 and vice versa at 3 and 5 days of lead time. For the CMIP6 evaluation purposes,
865 the models trained on ERA5 were slightly worse than the models trained on CMIP6. Surprisingly,
866 for evaluating on ERA5, models trained on CMIP6 were comparable, if not slightly better to the
867 ones trained on ERA5. These results are in line with results of [50], thus highlighting the dataset
868 usefulness of CMIP6 over ERA5. Note that the data’s raw size is roughly similar for both the datasets
869 as despite the ERA5’s temporal training range being 1979-2010 in this setup, it’s data availability
870 frequency is 1 hour compared to 6 hour in CMIP6.

871 To find out whether this superiority of CMIP6 over ERA5 is just a result of differences in temporal
 872 range, we conducted the similar study but with same dataset temporal characteristics (i.e. setting
 873 training years as 1979-2010 and subsampling the data at 6 hours). This time the results just for
 874 ResNet at 3 day lead time is shown in Table 2 and for all models at different lead times, is shown in
 875 Table 12. These results show that the performance is slightly worse for both the cases now. Thus
 876 showing that the performance improvement of training over CMIP6 than ERA5 is likely just the
 877 bigger temporal range.

878 D Visualizations

879 ClimateLearn provides visualization functionality to help with an intuitive understanding of model
 880 performance. Below is an example figure generated by ClimateLearn for visualizing the quality of
 881 a model’s forecast. Each row represents a distinct time in the test set. The leftmost column shows
 882 weather conditions at the time the model is making a prediction from. The next column shows the
 883 ground truth conditions at the forecast horizon. The next column shows the model’s predictions. The
 884 last column shows the model’s bias, and its per-pixel forecast error.

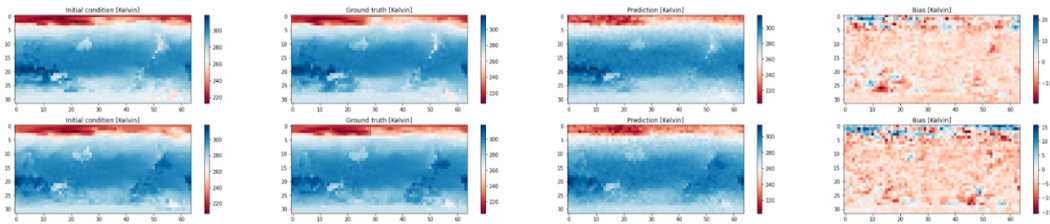


Figure 4: Example visualization of deterministic forecasting.

885 Additionally, ClimateLearn can generate the rank histogram for probabilistic forecasts. A rank
 886 histogram that resembles a uniform distribution means that the ground truth value is indistinguishable
 887 from any member of the forecast ensemble. A rank histogram that is skew right occurs when
 888 the ground truth is consistently lower than the ensemble prediction. A rank histogram that appears
 889 U-shaped is indicative of both low biases and high biases. An example figure generated by
 890 ClimateLearn for visualizing the rank histogram is shown below.

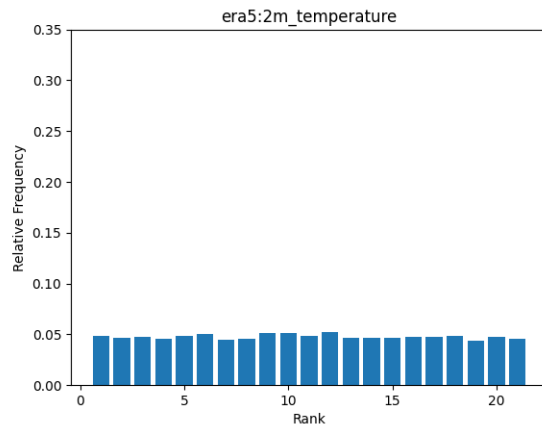


Figure 5: Example visualization of the rank histogram for probabilistic forecasting.

891 We show an example of how to generate another visualization called “mean-bias” in the next section.

892 E Code snippets

893 ClimateLearn can be used to download heterogeneous climate data from a variety of sources in a
894 single function call. Here, we provide an example for downloading ERA5 2-meter temperature data
895 at 5.625° resolution from WeatherBench.

```
1 from climate_learn.data import download
2 download(
3     root="./weatherbench-data",
4     source="weatherbench",
5     dataset="era5",
6     resolution="5.625",
7     variable="2m_temperature"
8 )
```

896 Further, ClimateLearn can process downloaded data into a form that is loadable into PyTorch.
897 In fewer than **30 lines**, the following code loads raw ERA5 data; normalizes it; splits it into train,
898 validation, testing sets; and prepares batches for the forecasting task.

```
1 # For flexibility with loading datasets and implementing new ones
2 # in the future, ClimateLearn's data processing pipeline is made
3 # up of three parts: the climate dataset (e.g., ERA5), the task
4 # (e.g., forecasting), and the PyTorch dataset (e.g., Map). These
5 # are all combined into a Pytorch Lightning DataModule.
6 from climate_learn.data.climate_dataset.args import ERA5Args
7 from climate_learn.data.task.args import ForecastingArgs
8 from climate_learn.data.dataset import MapDatasetArgs
9 from climate_learn.data import DataModule
10
11 # Next, we define the arguments: the location of the data, the
12 # variables we will use as input/target, and the data splits
13 root = "./weatherbench-data"
14 variables = ["2m_temperature"]
15 train_years = range(1979, 2016)
16 val_years = range(2016, 2017)
17 test_years = range(2017, 2019)
18
19 # Next, we construct the arguments for the three parts of the data
20 # processing pipeline to build the training dataset.
21 climate_dataset_args = ERA5Args(root, variables, train_years)
22 task_args = ForecastingArgs(
23     [f"era5:{var}" for var in variables], # format - dataset:var
24     [f"era5:{var}" for var in variables], # format - dataset:var
25     pred_range=72, # hours ahead to predict
26     history=3, # past time steps
27     subsample=6 # hours per time step
28 )
29 train_data_args = MapDatasetArgs(climate_dataset_args, task_args)
30
31 # The validation and test datasets can be constructed easily by
32 # copying arguments from the train dataset which are the same and
33 # modifying only what is needed.
```

```

34 val_data_args = train_data_args.create_copy({
35     "climate_dataset_args": {"years": val_years}
36 })
37 test_data_args = val_data_args.create_copy({
38     "climate_dataset_args": {"years": test_years}
39 })
40
41 # Finally, we can unify all parts of the data pipeline to get a
42 # single PyTorch Lightning data module.
43 dm = DataModule(train_data_args, val_data_args, test_data_args)

```

899 With the loaded data, ClimateLearn can be used to build, train, and evaluate a model in fewer than
900 **20 lines** of code.

```

1  import climate_learn as cl
2  from climate_learn.training import Trainer
3
4  model_kwargs = {
5      "in_channels": 1, # predicting 2m_temperature
6      "history": 3,    # matching 'ForecastingArgs'
7      "n_blocks": 4   # number of residual blocks to use
8  }
9  optim_kwargs = {} # use the default settings
10 mm = cl.load_forecasting_module(
11     data_module=dm,
12     model="resnet",
13     model_kwargs=model_kwargs,
14     optim_kwargs=optim_kwargs
15 )
16
17 trainer = Trainer()
18 trainer.fit(mm, dm)
19 trainer.test(mm, dm)
20

```

901 ClimateLearn can also be used to load pre-defined models (e.g., persistence, Rasp and Thuerey
902 [60]) as follows.

```

1  persistence = cl.load_forecasting_module(
2      data_module=dm,
3      preset="persistence"
4  )
5  rasp_theurey_2020 = cl.load_forecasting_module(
6      data_module=dm,
7      preset="rasp-theurey-2020"
8  )

```

903 ClimateLearn can use the trained forecasting models to produce visualizations in a single line of
904 code. For example, one visualization of interest is the mean bias, which shows the expected error of
905 the model's forecast, per pixel, over the evaluation period.

```
1 from climate_learn.utils.visualize import visualize_mean_bias
2 visualize_mean_bias(persistence, dm)
```

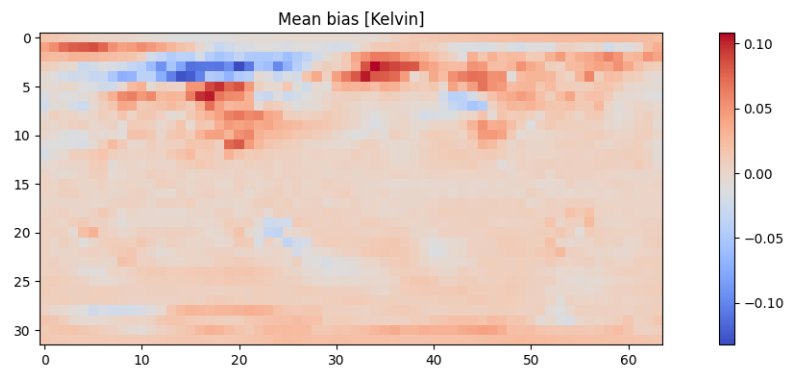


Figure 6: Visualization of the mean bias of temperature

906

907 This graphic shows that, on average, persistence has little bias below the equator. Over the northern
908 part of North America, persistence achieves negative mean bias, which means it generally under-
909 predicts 2-meter temperature in that region. Meanwhile, in the northern part of Europe, persistence
910 achieves positive mean bias, indicating overprediction.