# ClimaX:
# A foundation model for weather and climate

**Tung Nguyen**[1]**, Johannes Brandstetter**[2]**, Ashish Kapoor**[3]**,**
**Jayesh K. Gupta**[*2]**, and Aditya Grover**[*1]
[1]UCLA, [2]Microsoft, [3]Scaled Foundations

Most state-of-the-art approaches for weather and climate modeling are based on physics-informed numerical models of the atmosphere. These approaches aim to model the non-linear dynamics and complex interactions between multiple variables, which are challenging to approximate. Additionally, many such numerical models are computationally intensive, especially when modeling the atmospheric phenomenon at a fine-grained spatial and temporal resolution. Recent data-driven approaches based on machine learning instead aim to directly solve a downstream forecasting or projection task by learning a data-driven functional mapping using deep neural networks. However, these networks are trained using curated and homogeneous climate datasets for specific spatiotemporal tasks, and thus lack the generality of numerical models. We develop and demonstrate ClimaX, a flexible and generalizable deep learning model for weather and climate science that can be trained using heterogeneous datasets spanning different variables, spatio-temporal coverage, and physical groundings. ClimaX extends the Transformer architecture with novel encoding and aggregation blocks that allow effective use of available compute while maintaining general utility. ClimaX is pre-trained with a self-supervised learning objective on climate datasets derived from CMIP6. The pre-trained ClimaX can then be fine-tuned to address a breadth of climate and weather tasks, including those that involve atmospheric variables and spatio-temporal scales unseen during pretraining. Compared to existing data-driven baselines, we show that this generality in ClimaX results in superior performance on benchmarks for weather forecasting and climate projections, even when pretrained at lower resolutions and compute budgets. Source code is available at https://github.com/microsoft/ClimaX.
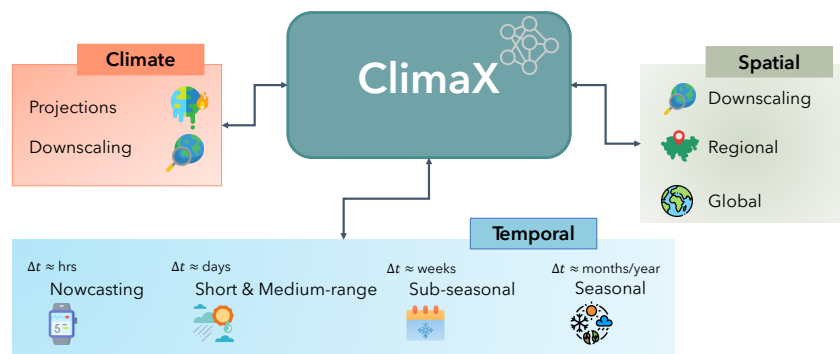
**Figure 1:** ClimaX is built as a foundation model for any weather and climate modeling task. On the weather front, these tasks include standard forecasting tasks for various lead-time horizons at various resolutions, both globally or regionally. On the climate front, making long term projections and obtaining downscaling results from lower resolution model outputs are standard tasks.

---

*Equal contributions as last authors, listed reverse alphabetically

Author email(s): tungnd@cs.ucla.edu, johannesb@microsoft.com, ashish.kapoor@gmail.com, jkg@cs.stanford.edu, adityag@cs.ucla.edu

# Contents

## 1. Introduction

Modeling weather and climate is an omnipresent challenge for science and society. With rising concerns around extreme weather events and climate change, there is a growing need for both improved weather forecasts for disaster mitigation and climate projections for long-term policy making and adaptation efforts [MD+21]. Currently, numerical methods for global modeling of weather and climate are parameterized via various general circulation models (GCM) [Lyn08]. GCMs represent system of differential equations relating the flow of energy and matter in the atmosphere, land, and ocean that can be integrated over time to obtain forecasts for relevant atmospheric variables [Lyn08; BTB15]. While extremely useful in practice, GCMs also suffer from many challenges, such as accurately representing physical processes and initial conditions at fine resolutions, as well as technological challenges in large-scale data assimilation and computational simulations [Bau+20]. These factors limit their use in many scenarios, especially in simulating atmospheric variables quickly at very short time scales (e.g., a few hours) or accurately at long time scales (e.g., beyond 5-7 days) [Zha+19].

In contrast, there has been a steady rise in data-driven approaches for forecasting of atmospheric variables, especially for meteorological applications [GKH15; DB18; Web+20; SM19; Sch18; Kas+21; Sch+21; Rei+19; Hun+19; Sch+17]. The key idea here is to train deep neural networks to predict the target atmospheric variables using decades of historical global datasets, such as the ERA-5 reanalysis dataset [Her+20]. Unlike GCMs, these networks are not explicitly grounded in physics, and lack general-purpose utility for Earth system sciences as they are trained for a specific predictive modeling task. Yet, with growing compute and datasets, there is emerging evidence that these models can achieve accuracies competitive with state-of-the-art numerical models in many scenarios, such as nowcasting of precipitation [Rav+21; Søn+20] and medium-range forecasting of variables like temperature, wind and humidity [WDC20; RT21; Kei22; Pat+22; Bi+22; Lam+22]. While these trends are encouraging, there remain concerns regarding the generality of such data-driven methods to diverse real-world scenarios, such as forecasting of extreme weather events and longer-term climate projections, especially under limited spatiotemporal supervision and computational budgets.

Variants of the aforementioned challenges apply broadly throughout machine learning (ML). In disciplines such as natural language processing and computer vision, it is well acknowledged that ML models trained to solve a single task using supervised learning are label-hungry during training and brittle when deployed outside their training distribution [Tao+20]. Recent works have shown that it is possible to mitigate the supervision bottleneck by *pretraining* [Dev+18; He+22] large unsupervised "foundation" models [Bom+21] on huge passive datasets, such as text and images scraped from the internet [Ram+22; Bro+20; Liu+21; Ree+22b]. Post pretraining, there are many ways to *finetune* the same model on arbitrary target task(s) with little to none (i.e., zero-shot) additional supervision. Besides low target supervision, these models also generalize better to shifts outside their training distribution [Hen+20a; Zha+22b], improving their reliability.

Inspired by the above successes, this work studies the question: how do we design and train a foundation model for weather and climate that can be efficiently adapted for general-purpose tasks concerning the Earth's atmosphere? We propose ClimaX, a foundation model for weather and climate. For pretraining any foundation model, the key recipe is to train a deep architecture on a large dataset using an unsupervised objective. For example, many foundation models for language and vision train large transformers on Internet-scale datasets using generative modeling. While conceptually simple, this scaling recipe is riddled with challenges for weather and climate domains, that we discuss below and propose to resolve with ClimaX.

First, it is unclear what constitutes an Internet-scale passive dataset for pretraining ClimaX. The size of historical weather and climate datasets at any given time is fixed and increases at an almost constant rate everyday, as it corresponds to processed sensor measurements of naturally occurring phenomena. Our first key proposal is to go beyond these datasets to explicitly utilize physics-informed climate simulation models. Many such models are in use today, for example, the CMIP6 collection [Eyr+16] of climate modeling simulations consists of runs of ∼100 distinct climate models from 49 different climate modeling groups. We show that the heterogeneity in these simulation datasets serves as a source of rich and plentiful data for pretraining ClimaX.

Second, we need a model architecture that can aptly embrace the heterogeneity of the above climate datasets. Climate data is highly multimodal, as observations typically correspond to many different, unbounded variables

with varying datatypes (e.g., pressure, temperature, humidity). Moreover, many observational datasets are irregular in the sense that they differ in their spatiotemporal coverage and might correspond to different subsets of atmospheric variables. We resolve the above challenges in ClimaX by repurposing the vision transformer [Dos+20; Vas+17]. In contrast to earlier work where the input data is represented as an image with different atmospheric variables treated as the channels thereof [Pat+22; Bi+22], we treat them as separate modalities to enable more flexible training even with irregular datasets. This has the side-effect of drastically increasing the sequence length, which we propose to resolve via a cross-attention style channel aggregation scheme prior to the self-attention layers.

Third and last, we need a pretraining objective that can learn complex relationships between the atmospheric variables and permit effective finetuning for downstream tasks. Given the spatiotemporal nature of climate data, we propose a randomized forecasting objective for pretraining ClimaX. Here, the goal of the model is to forecast an arbitrary set of input variables at an arbitrary time into the future. While simple and intuitive, we show that such a pretraining objective aids finetuning to novel tasks and timescales even beyond the pretraining window, such as sub-seasonal to seasonal cumulative predictions, climate projections, and downscaling of climate models. See Figure 1 for a list of tasks considered in this work.

Empirically, we demonstrate that a single pretrained model can be finetuned for many tasks (e.g., multi-scale weather forecasting, climate projections, downscaling) under a range of operating conditions involving different spatiotemporal resolutions, geographical regions, and target prediction variables, including those unseen during training. Notably, our benchmark results are state-of-the-art on ClimateBench [WP+22] and competitive with the operational Integrated Forecasting System (IFS) [Wed+15] on WeatherBench [Ras+20], even when our model is trained on moderate resolutions using only a maximum of 80 NVIDIA V100 GPUs.

Finally, we show promising scaling laws of ClimaX with natural axes of performance improvements for larger number of pre-training datasets, larger models, and scaling to higher resolution gridded datasets. While especially the last is in line with recent and concurrent works on data-driven weather forecasting [Pat+22; Bi+22; Lam+22], to the best of our knowledge, ClimaX is the first of its kind data-driven model that can effectively scale using heterogeneous climate datasets during pretraining, and generalize to diverse downstream tasks during finteuning, paving the way for a new generation of data-driven models for Earth systems science.

## 2. Background and Related Work

Current weather and climate models in use today rely extensively on numerical methods and computational simulations to predict and understand the Earth's weather and climate systems. These tasks include various *numerical weather prediction* (NWP) systems which use computer simulations to make short-term forecasts of weather conditions as well as climate models which use similar techniques to simulate and predict the long-term changes in the Earth's climate. Most notably, at the core of both weather and climate models lie the same set of primitive equations.

For climate modeling, earth system models (ESM) [Hur+13], or "coupled models", that couple together simulations which govern the atmosphere, cryosphere, land, and ocean processes are considered the state-of-the-art. Primarily these simulations are based on general circulation models (GCMs) [Sat04; Lyn08; Ado14; MD+21] which date back to the works of Phillips [Phi56] and Lorenz [Lor67] solving Navier-Stokes equations on a rotation sphere to model fluid circulation. These models are often used to perform various *factor sensitivity* studies to examine how the changes in certain forcing factors like greenhouse gas concentrations can affect the global or regional climate and help in *climate projections* to help understand future conditions.

Numerical Weather Prediction (NWP) models share many components of GCMs, especially the atmospheric components [BTB15; Lyn08; Kal03]. However, incorporating *data assimilation* [LSZ15; Gro22] which involves combining observations and various measurements of the atmosphere and oceans together with these numerical models is important for accurate forecasts and simulations. Another significant distinction between weather and climate models is the framing of the solution for underlying equations: *initial value problem* for weather, while *boundary value problem* for climate [BTB15]. Different difficulty levels of these solution approaches results

in the fact where climate models tend to be global often at coarser spatio-temporal resolutions while weather models can range from global to local and regional models of very high spatio-temporal resolutions [War10].

Despite their noted success, including the recent 2021 Nobel Prize in Physics [RRH22], there is considerable debate around the limitations of general circulation models (GCMs), particularly structural errors across models and the fact that current GCMs are designed to reproduce observed climate [Bal+22]. The climate science community has been aware of these challenges which resulted in the creation of Coupled Model Intercomparison Project (CMIP) as a standardized protocol for evaluating and comparing the performance of different climate models [Mee+00]. As we will see in the following sections, not only has CMIP been playing a crucial role in the advancement of our understanding of climate change and its potential impacts, its evaluation procedure has resulted in enormous quantity of data making modern deep learning based approaches quite attractive for many tasks. Notably, encoding this knowledge into a "foundation" machine learning model with much faster inference and data assimilation capabilities can pave the way for a much wider impact.

## 2.1. Data sources

Unlike data in computer vision or natural language processing, weather and climate data is not solely based on sensed data, instead incorporates information from a diverse range of sources. For example, *reanalysis* weather data blends meteorological observations with past short-range weather forecasts via data assimilation [BTB15]. The data measurements themselves are highly heterogeneous, representing various physical variables with different data types (e.g. pressure, temperature, humidity) that are recorded at different, relatively sparse, spatial locations at different temporal frequencies. These measurements can be integrated together with known physics inform the design of climate simulations, which again produce data with different variables at different scales. From a machine learning perspective, the plethora of available data thus spans multiple axes: from direct weather measurements at land, sea, or atmosphere, over multiple decades of re-analyzed weather data at different spatial scales, to physics-informed climate projections for various scenarios. Most notably, the data shares the same set of primitive equations, but with fairly different characteristics. Below we describe two of the most commonly used data sources for weather and climate modeling.

### 2.1.1. CMIP6

The Coupled Model Intercomparison Project (CMIP) [Mee+00] is an international effort across different individual climate modeling groups to come together to compare and evaluate their global climate models. While the main goal of CMIP is to improve the understanding of Earth's climate system and improve the accuracy of its simulations, the recent data from their experimental runs is easily accessible on the CMIP6 [Eyr+16] archive. In CMIP6, where "6" refers to the most recent phase of the project, 49 groups are involved with their experiments covering wide range of climate variables including temperature, precipitation, sea level and others from hundreds of models. This results in global projections of various climate scenarios from as early as 1850 onwards, all following similar governing equations, but with different *forcings*, e.g., greenhouse gas emissions that affect the climate.

### 2.1.2. ERA5

The ERA5 reanalysis archive [Her+18; Her+20] of the European Center for Medium-Range Weather Forecasting (ECMWF) is the predominant data source for learning and benchmarking weather forecasting systems. Once completed, the ERA5 reanalysis is set to embody a detailed record of the global atmosphere, land surface and ocean waves from 1950 onwards. The currently available ERA5 reanalysis data combines the state of the art forecasting model called Integrated Forecasting System (IFS) [Wed+15] of ECMWF with available observations to provide the best guess of the state of the atmosphere, ocean-wave and land-surface quantities at any point in time. In its raw form, the available reanalyzed data is huge: 40 years, from 1979 to 2018, on a $0.25° \times 0.25°$ global latitude-longitude grid of the Earth's sphere, at hourly intervals with different climate variables at 37 different altitude levels plus the Earth's surface. The grid overall contains $721 \times 1440$ grid points for latitude and longitude, respectively. The altitude levels are presented as pressure levels.

## 2.2. Tasks

Given the scale of data availability, increasing compute requirements of current numerical methods despite it being difficult to incorporate real observational data into them, machine learning is increasingly finding applications in many of the tasks related to weather and climate modeling. When it comes to **weather**, the main task of interest is *forecasting* the future values of key weather variables. These tasks can take the following forms depending on temporal and spatial horizons of interest:

- **Global forecasting** tasks that range from a few hours (i.e., nowcasting) to days and weeks in lead time (i.e., short and medium range forecasting). Often these tasks are evaluated on the ERA5 reanalysis dataset (see Section 2.1.2) with Operational IFS [Wed+15] of the European Center for Medium-Range Weather Forecasting (ECMWF) being the current state-of-the-art NWP baselines.
- **Regional forecasting** tasks which could range from weather forecasting in continental North America or Europe to individual state, county or city.
- **Sub-seasonal to seasonal prediction (S2S)** [VR18; Vit+22] which is the task of forecasting the weather with lead times between 2 weeks and 2 months. S2S bridges the gap between weather forecasting and seasonal climate prediction, and is critical to disaster mitigation. Often at such long horizons, predicting instantaneous values of key weather variables can be a difficult task and therefore the focus is often on averaged value of key weather variables over a certain time horizon, e.g. weekly average precipitation.

Whereas deep learning approaches for regional or S2S tasks are scarce, most of the recent and concurrent work focuses on global forecasting tasks. Rasp and Thuerey [RT21] were the first to use pretraining on climate simulations to achieve good data-driven medium-range weather prediction with a ResNet [He+16], Weyn, Durran, et al. [WDC20] used CNNs on a cubed sphere for global weather prediction, Weyn, Durran, et al. [Wey+21] forecast weather sub-seasonally with a large ensemble of deep-learning weather prediction models, Keisler [Kei22] applied a graph neural network based approach to weather forecasting, Ravuri, Lenc, et al. [Rav+21] use deep generative models of radar for precipitation nowcasting, Arcomano, Szunyogh, et al. [Arc+20] build a reservoir computing-based, low-resolution, global prediction model, and MetNet [Søn+20] takes as input radar and satellite data to forecast probabilistic precipitation maps. These approaches are complemented by general machine learning models for fluid dynamics [Li+20; Koc+21; Lu+21; Bra+22; BWW22]. Finally, recent state-of-the-art neural weather models such as FourCastNet [Pat+22], Pangu-weather [Bi+22], or GraphCast [Lam+22], which also perform global forecasting tasks, use the highest resolution 0.25° ERA5 data, and are optimized on the respective hardware resources.

On the other hand, **climate** tasks have to deal with much longer time horizons. Possible categories of tasks where machine learning can help include climate projection and climate model downscaling:

- **Climate projection** is the task of generating estimates of climate change under different future socio-economic scenarios. Usually, this takes the form of figuring out the response of the climate system to different forcing factors such as greenhouse gases and aerosol emissions. Climate projection is a crucial task in understanding and preparing for the potential impacts of climate change.
  While the application of machine learning in this field is still in its early stages, recent efforts have been made to standardize evaluation in this domain. One example of this is ClimateBench [WP+22], which is a benchmark dataset drawing on CMIP6 to provide an evaluation framework for machine learning models that aim to improve the accuracy of climate projections. This benchmark aims to provide a consistent and reliable evaluation method for various machine learning models that are applied to climate projections.
- A more popular application of ideas in machine learning is towards **downscaling** of climate model. Global climate models typically have a coarse spatial resolution, which means that they can only provide a rough estimate of climate conditions at a local or regional scale. Moreover, the simulations often reflect systematic biases that deviate from trends in the observation data. The aim of climate model downscaling is to create locally accurate climate information from global climate projections by relating those to observed local climatological conditions. This process improves the spatial and temporal resolution of the data, making it more suitable for use in local and regional analyses. Downscaling methods can be divided into

*dynamic* approaches that relate outputs of global climate models with those of regional climate models, and *statistical* approaches that infer the desired transformations using data-driven approaches [WW97]. Dynamic approaches are physically consistent, but can be slow and have large biases, whereas statistical approaches need large amounts of data to learn expressive mappings that hold for target output scenarios.

Similar to weather forecasting, deep learning has emerged as appealing alternative in climate science as well. Recent approaches comprise surrogate models to emulate climate projections [Web+20; SM19; Sch18; BGS20; Man+20], extract contextual cues from existing datasets or simulations [Rei+19; Hun+19; Sch+17], and perform climate model downscaling [Sac+18; Van+17; BMMG20]. Climate model downscaling usually inputs low-resolution reanalysis data and local orographic information to obtain high-resolution local information. Many recent approaches are based on convolutional architectures [Höh+20; Vau+21; Mar+22].

## 2.3. Foundation models

Bommasani, Hudson, et al. [Bom+21] gave the term "foundation models" to the emerging paradigm of training scalable deep learning models on broad data via self-supervision which could then be adapted (often via finetuning) to a wide range of downstream tasks. Current notable examples include BERT [Dev+18], GPT [Bro+20] and PaLM [Cho+22], in language, CLIP [Rad+21], Florence [Yua+21], BEiT [Wan+22] for vision-language. Outside applications on data crawled from web, this paradigm has also started finding success in various scientific domains like protein design [Ver+22]. Key significance of such models has been identified as *emergence* with respect to model capabilities and *homogenization* with respect to methodologies for different tasks, domains, and modalities, enabled by the principles of transfer learning [TP12] at scale. While a foundation model itself should be considered incomplete, it can provide a common basis from which various task-specific models can be derived. Current research at the intersection of weather and climate science and ML has largely focused on designing separate models for every task of interest despite potential availability of fairly diverse large scale data with shared underlying physics and geology across these tasks. A few recent works have proposed pretraining techniques for satellite imagery and remote sensing [YL20; Con+22; Ree+22a] but they have so far not been applied to multi-sensory data and variables in weather and climate.

## 3. Approach

Given the availability of large scale data sources, together with shared physics and geology between various weather and climate tasks, we aim to build a generalizable deep learning foundation model. The model needs to be able to input heterogeneous datasets of different variables, and provide spatio-temporal coverage based on physical groundings. We, therefore, first take a closer look at input representations, and next design a model to cope with their heterogeneity - local, global, and across variables.

### 3.1. Input representation

We are interested in gridded prediction tasks, in which the model takes an input of shape $V \times H \times W$ and predicts an output of shape $V' \times H' \times W'$. $V$ refers to the number of input variables, which can be weather conditions such as geopotential and temperature, or climate forcing factors such as $CO_2$ and $SO_2$. $H$ and $W$ refer to the spatial resolution of the input data, which depends on how densely we grid the globe. This general representation captures a broad variety of downstream tasks in Earth systems science. Similarly, $V', H', W'$ refer to the variables and spatial resolution of the predicted outputs. We mainly work with two spatial resolutions: 5.625° ($32 \times 64$ grid points) and 1.40625° ($128 \times 256$ grid points). Semantically, a $H \times W$ map can represent the entire globe or a specific region such as North America.

### 3.2. Model architecture

We aim to design a foundation model that we can pretrain on heterogeneous data sources and then finetune to solve various downstream weather and climate tasks. From Section 3.1, one could think of the tasks as
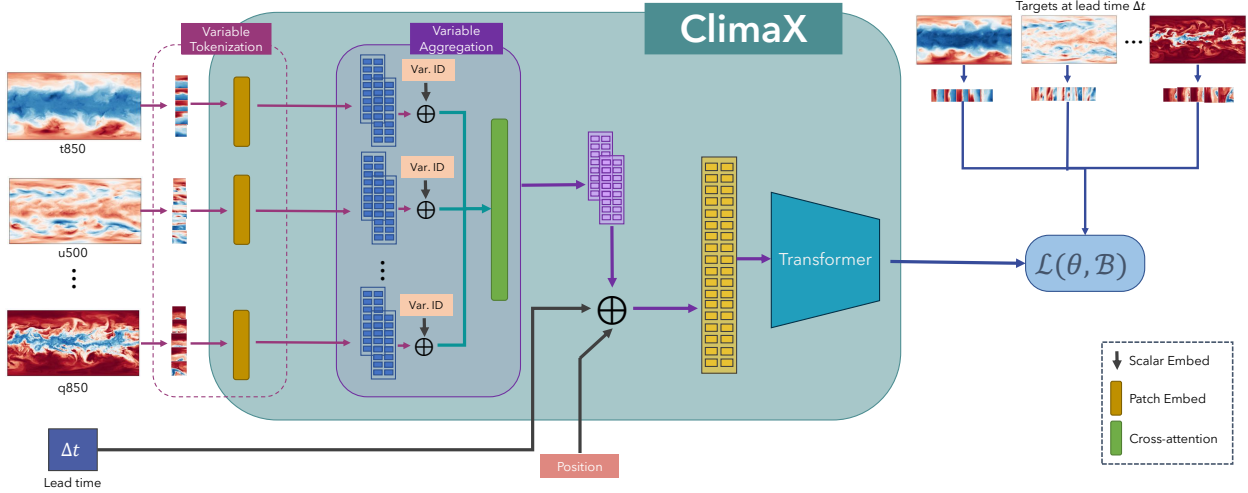
**Figure 2:** Pretraining phase of ClimaX. Variables are encoded using variable-separate tokenization, and subsequently aggregated using variable aggregation. Together with position embedding and lead time embedding those are fed to the ViT backbone.

image-to-image translation problems with $V$ input channels and $V'$ output channels. This makes any image architecture a natural fit, such as UNet [RFB15], ResNet [He+16], or Vision Transformers (ViT) [Dos+20]. However, the settings of climate and weather tasks are much broader, where we may want to make predictions for regional or even spatially incomplete data, forecast unseen climate variables, or finetune the model on data at different resolutions from pretraining. Current CNN-based architectures are not applicable in these scenarios, as they require the input to be perfectly gridded, contain a fixed set of variables, and have a fixed spatial resolution. Transformers-based architectures, on the other hand, provide much better flexibility by treating the image-like data as a set of tokens. Therefore, we build ClimaX architecture upon Vision Transformers (ViT) [Dos+20; Vas+17], and propose two major architectural changes, namely *variable tokenization* and *variable aggregation* to further improve the flexibility and generality, which we will describe next.

### 3.2.1. Variable tokenization

Given an input of shape $V \times H \times W$, ViT tokenizes the input into a sequence of $(H/p) \times (W/p) = h \times w$ patches, with each patch having a size of $V \times p^2$, where $p$ is the patch size. This tokenization scheme works well for image data, as $V$ is always the RGB channels, which is the same for all datasets. However, this is not true for climate and weather data, where the number of physical variables can vary between different datasets. For example, in the CMIP6 project [Eyr+16], each dataset contains simulated data of a different climate model, and thus has a different set of underlying variables. Therefore, we propose *variable tokenization*, a novel tokenization scheme that tokenizes each variable in the input separately. Specifically, each input variable as a spatial map of shape $H \times W$ is tokenized into a sequence of $h \times w$ patches, which results in $V \times h \times w$ patches in total. Finally, each input patch of size $p^2$ is linearly embedded to a vector of dimension $D$, where $D$ is the chosen embedding size. The output of the variable tokenization module therefore has a dimension of $V \times h \times w \times D$. Figure 3 illustrates our proposed tokenization scheme.

### 3.2.2. Variable aggregation

While variable tokenization allows ClimaX to learn from datasets with varying numbers of input variables, it has two inherent problems. First, it results in a sequence of length $V \times h \times w$ which increases linearly with the number of variables. Since we use attention to model the sequence, the memory complexity scales quadratically with the number of variables. This is computationally expensive, as we can have up to 48 input
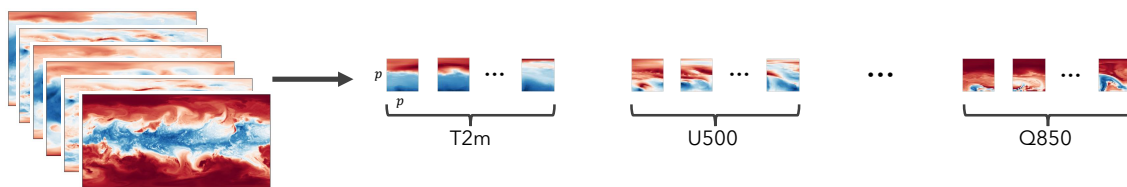
**Figure 3:** Variable tokenization. Each variable is independently tokenized.

variables in our experiments. Moreover, because we tokenize each variable separately, the input sequence will contain tokens of different variables with very different physical groundings, which can create difficulties for the attention layers to learn from. We therefore propose *variable aggregation* to solve the two mentioned challenges. For each spatial position in the $h \times w$ map, we perform a cross-attention operation, in which the query is a learnable vector, and the keys and values are the $V$ embedding vectors of $V$ variables at that position. The cross-attention module outputs a single vector for each spatial position, thus reducing the sequence length to $h \times w$, significantly lowering the computational cost. Moreover, the sequence now contains unified tokens with universal semantics, creating an easier task for the attention layers. Figure 4 shows our proposed variable aggregation.
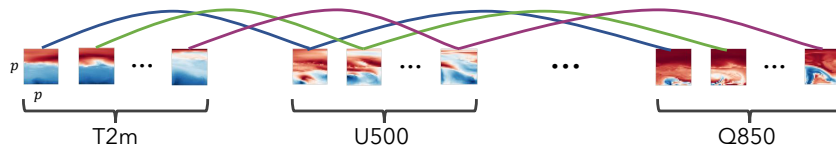


**Figure 4:** Position-based variable aggregation reduces a sequence of length $V \times h \times w$ to $h \times w$.

### 3.2.3. Transformer

Post variable aggregation, we need a sequence model for generating the output tokens. While in principle, one could use any general sequence model, we propose to extend a standard Vision Transformer (ViT). Moreover, since the standard ViT treats image modeling as pure sequence-to-sequence problems, it can perform tasks that some other variations cannot [Liu+21; Liu+22], such as learning from spatially incomplete data, where the input does not necessarily form a complete grid. This is useful in the regional forecasting task we consider in Section 4.2.2. In the experiments, we report results with 8 attention layers, an embedding size of 1024, and a hidden dimension of $1024 \times 4$. After the attention layers, we employ a prediction head that takes a token and outputs a vector of size $V' \times p^2$. The prediction head is a 2-layer MLP with a hidden dimension of 1024. We provide more details in Appendix A.

### 3.3. Datasets

### 3.3.1. Pretraining

We believe that CMIP6's diversity and scale presents an attractive opportunity for pretraining large-scale foundation models. However, handling the inconsistent set of variables across different data sources can be a challenge. In this work we only use a subset of variables from five different data sources (MPI-ESM, TaiESM, AWI-ESM, HAMMOZ, CMCC) containing global projections of climate scenarios from 1850 to 2015 with the time delta of 6 hours as described in Table 8. Due to variable original resolution, we choose to simplify our data-loading by regridding them to commonly used resolutions [Ras+20; RT21] of 5.625° ($32 \times 64$ grid points) and 1.40625° ($128 \times 256$ grid points)[1].

---

[1]Regridding was done using the xesmf Python package [Zhu18] using bilinear interpolation.

### 3.3.2. Finetuning and evaluation

We use the ERA5 reanalysis data as described in Appendix C.2, as the source of datasets for finetuning and evaluation for various weather related downstream tasks. Due to its large size, it is common to regrid [Ras+20; RT21] the high-resolution data to lower resolutions like 5.625° ($32 \times 64$ grid points) and 1.40625° ($128 \times 256$ grid points) to fit within the available computational constraints[2]. We follow the evaluation procedure by Rasp and Thuerey [RT21] and use this data to assess the forecasting performance of our ML models at different lead time horizons. More details about the individual datasets are in their appropriate experiment sections.

### 3.4. Training

### 3.4.1. Pretraining

We pretrain ClimaX on CMIP6 data to predict future weather conditions given the current conditions. That is, given the weather snapshot $X_t$ of shape $V \times H \times W$ at a particular time $t$, ClimaX learns to predict the future weather scenario $X_{t+\Delta t}$ of the same shape at lead time $\Delta t$. To obtain a pretrained model that is generally applicable to various temporal forecasting tasks, we randomize the lead time from 6 hours to 168 hours (i.e., 1 week) during pretraining. We add the lead time embedding to the tokens to inform the model of how long it is forecasting into the future. The lead time embedding module is a single-layer MLP that maps a scalar to a vector of the embedding size $D$. Figure 2 depicts the forward pass of ClimaX in pretraining. For an input $X_t$, we sample a lead time $\Delta t \sim \mathcal{U}[6, 168]$ and get the corresponding ground truth $X_{t+\Delta t}$. Input variables are tokenized separately using variable tokenization, and are subsequently aggregated at each spatial location, resulting in a sequence of $h \times w$ unified tokens. We add the tokens with the lead time embedding and positional embedding before feeding the sequence to the ViT backbone. The output of the last attention layer is fed to a prediction head, which transforms the sequence back to the original shape of $V \times H \times W$.

We employ the latitude-weighted mean squared error [Ras+20] as our objective function. Given the prediction $\tilde{X}_{t+\Delta t}$ and the ground truth $X_{t+\Delta t}$, the loss is computed as:

$$\mathcal{L} = \frac{1}{V \times H \times W} \sum_{v=1}^{V} \sum_{i=1}^{H} \sum_{j=1}^{W} L(i)(\tilde{X}_{t+\Delta t}^{v,i,j} - X_{t+\Delta t}^{v,i,j})^2, \tag{1}$$

in which $L(i)$ is the latitude weighting factor:

$$L(i) = \frac{\cos(\mathrm{lat}(i))}{\frac{1}{H} \sum_{i'=1}^{H} \cos(\mathrm{lat}(i'))}, \tag{2}$$

where $\mathrm{lat}(i)$ is the latitude of the corresponding $i$th row of the grid. The latitude weighting term accounts for the non-uniformity in areas when we grid the round globe. Grid cells toward the equator have larger areas than the cells near the pole, and thus should be assigned more weights.

### 3.4.2. Finetuning

ClimaX has four learnable components, including the token embedding layers, the variable aggregation module, the attention blocks, and the prediction head. We evaluate the performance of ClimaX on various downstream tasks, which we categorize into two finetuning scenarios: one in which the downstream variables belong to the set of pretraining variables, and the other with variables unseen during pretraining. In the first case, we finetune the entire model, and in the latter, we replace the embedding layers and the prediction head with newly initialized networks, and either finetune or freeze the other two components. We present more details of each downstream task in Section 4.

---

[2]Regridding was done using the xesmf Python package [Zhu18] using bilinear interpolation.

## 4. Experiments

We finetune ClimaX on a diverse set of downstream tasks to evaluate its performance and generality. We categorize the tasks into forecasting, climate projection, and climate downscaling. The experiments aim to answer the following questions:

- How does ClimaX perform on global forecasting compared to the current state-of-the-art NWP system?
- Can we finetune ClimaX to make forecasts for a specific region or at different temporal horizons from pretraining?
- How well does ClimaX perform on climate tasks that are completely different from pretraining?

In addition to the main experiments, we analyze the scaling property of ClimaX, i.e., how the performance of ClimaX improves with increasing data size, model capacity, and data resolution. Finally, we perform comprehensive ablation studies to understand the trade-off between computation and performance when finetuning ClimaX.

### 4.1. Neural baselines

In global forecasting, we compare ClimaX with IFS [Wed+15], the current gold standard in weather forecasting. In tasks we do not have a baseline, we compare with UNet [RFB15; GB22] and ResNet [He+16], two CNN baselines commonly used in vision tasks. We borrow the ResNet architecture from Weatherbench [Ras+20]. The exact architectural details of these baselines are in Appendix A.2.

### 4.2. Forecasting

### 4.2.1. Global forecasting

Given global weather conditions $X_t$ at a particular time $t$, we want to forecast the weather at a future time $X_{t+\Delta t}$, in which $\Delta t$ is the lead time. The input variables include 6 atmospheric variables at 7 vertical levels, 3 surface variables, and 3 constant fields, resulting in 48 input variables in total. The details of the variables are in Table 9. We evaluate ClimaX on predicting four target variables: geopotential at 500hPa (Z500), the temperature at 850hPa (T850), the temperature at 2 meters from the ground (T2m), and zonal wind speed at 10 meters from the ground (U10). Z500 and T850 are the two standard verification variables for most medium-range NWP models and are often used for benchmarking in previous deep learning works, while the two surface variables, T2m and U10, are relevant to human activities. We consider seven lead times: 6 hours, $\{1, 3, 5, 7\}$ days, 2 weeks, and 1 month, which range from nowcasting to short and medium-range forecasting and beyond. We consider predicting each target variable at each lead time a separate task, and finetune a separate model for each task. We discuss alternative finetuning protocols in Section 4.6.

We compare ClimaX with IFS and the two CNN baselines on the ERA5 dataset at both 5.625° and 1.40625° resolutions. Following [Ras+20], we split the data into three sets, in which the training data is from 1979 to 2015, the validation data is in 2016, and the test data is in 2017 and 2018. We finetune ClimaX and train the other deep learning baselines using the latitude-weighted MSE loss in Equation (1). We perform early stopping on the validation loss for all deep learning models, and evaluate the best checkpoint on the test set. For IFS, we download the predictions from the TIGGE archive [Bou+10] for the year 2018[3]. We compare all methods on latitude-weighted root mean squared error (RMSE) and latitude-weighted anomaly correlation coefficient (ACC), two commonly used metrics in previous works. The formulations of the two metrics are in Appendix D.1. Lower RMSE and higher ACC indicates better performance.

Figures 5 and 6 show the performance of ClimaX and the baselines at 5.625° and 1.40625°, respectively. At low resolution, IFS outperforms ClimaX on 6-hour to 5-day prediction tasks. On longer horizons, however, ClimaX performs comparably to or slightly better than IFS, especially on 14-day prediction. At higher resolution, the

---

[3]We were not able to download IFS predictions for 2017 due to some server issues.
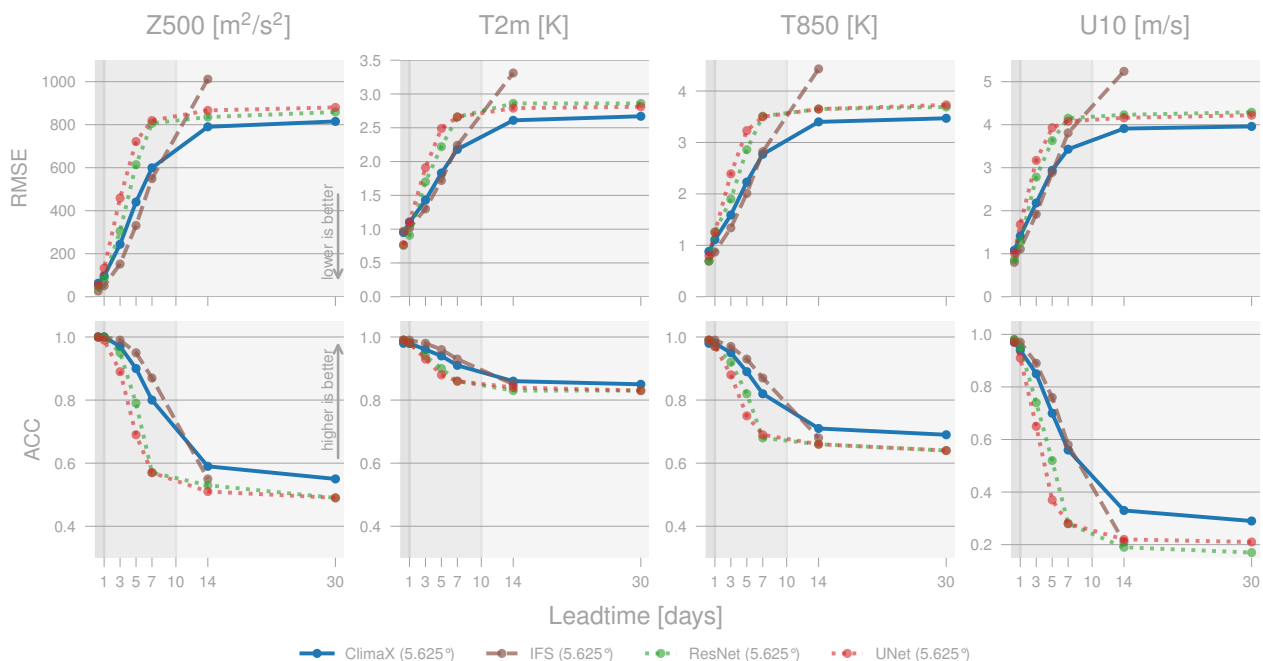
**Figure 5:** Performance on global forecasting on ERA5 at 5.625°.

performance of ClimaX closely matches that of IFS even for short horizons, and is superior in forecasting at 7 days and beyond. The trends are similar for both RMSE and ACC. The two CNN baselines perform similarly and achieve reasonable performance, but lag behind ClimaX and IFS on all tasks. We include other additional task-specific baselines [Pat+22; Bi+22; Lam+22] in Appendix D.2. These baselines are trained on higher-resolution ERA5 (0.25°) so are not directly comparable.

### 4.2.2. Regional forecasting

It is not always possible to make global predictions, especially when we only have access to regional data In this section, we evaluate ClimaX on *regional forecasting* of the relevant variables in North America, where the task is to forecast the future weather in North America given the current weather condition in the same region. We create a new dataset from the ERA5 data at 1.40625° that has the same set of variables but just focuses on the North America region. We call this dataset ERA5-NA and present details of how to construct it in Appendix C.2. Training, validation, and test splits are done similarly to Section 4.2.1. Figure 7 illustrates the finetuning process of ClimaX on this task, where the only difference from global forecasting is the input now only contains tokens that belong to North America.

Since the task has not been considered in previous works, we compare ClimaX with the two CNN baselines ResNet and UNet, and the scratch-trained version of ClimaX, which we refer to as Cli-ViT. In addition, we finetune two ClimaX models, in which one was pretrained on CMIP6 at 1.40625°, and the other was pretrained on 5.625° data. To finetune the low-resolution model on higher-resolution data, we follow the common practice of interpolating the positional embedding [Dos+20; Tou+21]. We denote this model as ClimaX-pos-interp. We evaluate all methods on predicting Z500, T2m, and T850 at lead times of 3, 5, and 7 days. Latitude-weighted RMSE is used as the evaluation metric.

Figure 8 compares the performance of ClimaX and the baselines. ClimaX is the best performing method among different target variables and lead times. Interestingly, even though pretrained on data at a lower resolution, ClimaX-pos-interp achieves the second best performance in predicting Z500 and T850, and only underperforms ResNet in predicting T2m at 3-day lead time. This result shows that ClimaX can gain strong
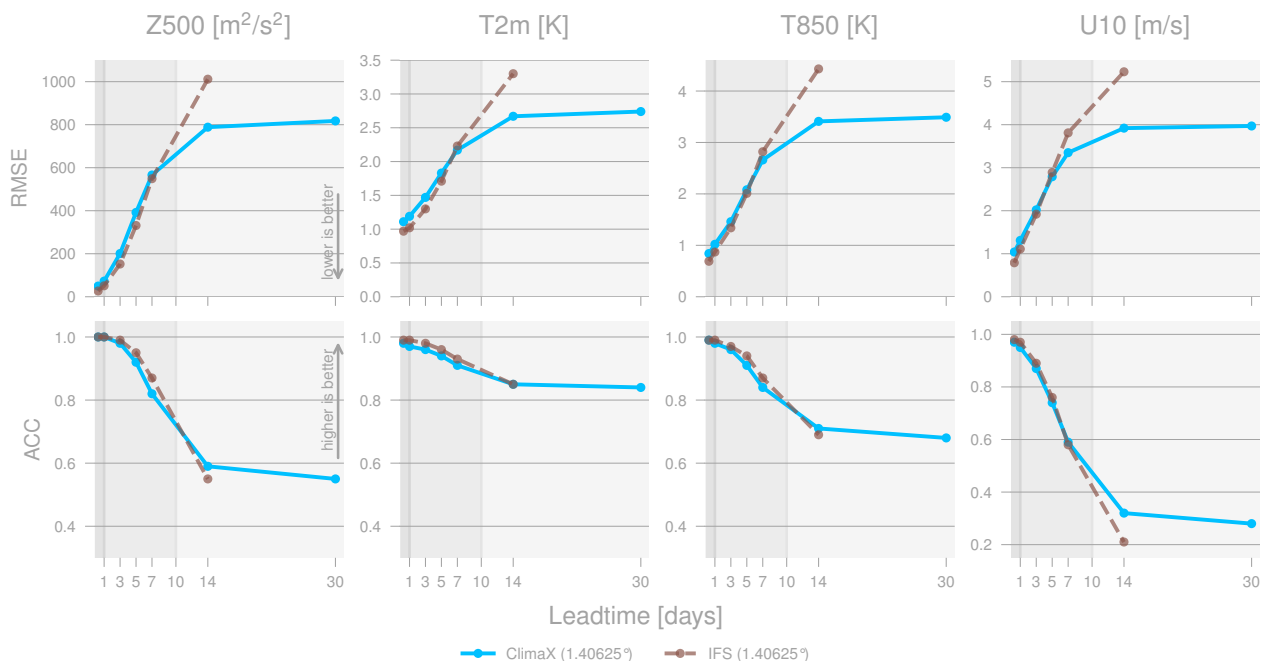
**Figure 6:** Performance on global forecasting on ERA5 at 1.40625°.

performance on tasks that have different spatial coverage or even different spatial resolution from pretraining.

### 4.2.3. Sub-seasonal to seasonal cumulative prediction

Sub-seasonal to seasonal (S2S) prediction is the task of forecasting at a time range between 2 weeks and 2 months [VR18], which bridges the gap between weather forecasting and climate projection. Compared to the other two well-established tasks, S2S prediction has received much less attention, despite having a significant socioeconomic value in disaster mitigation. Recent works have proposed data-driven approaches based on both traditional machine learning [Hwa+19; Pro+18; TL18] and deep learning [Wey+21; Zho+21; Ore+19], but their performances often lag behind adaptive bias correction methods [Mou+23a] on standard benchmarks [Mou+23b]. Here, following the S2S competition (https://s2s-ai-challenge.github.io/), we aim to predict the biweekly average statistics of weeks 3-4 and weeks 5-6, which correspond to lead times of 2 weeks and 4 weeks, respectively. We construct ERA5-S2S, a new dataset from 5.625° ERA5 that has the same input variables, but the output variables are averaged from the lead time to 2 weeks ahead into the future.

We compare ClimaX with ResNet, UNet, and Cli-ViT on the S2S prediction of four target variables: T850, T2m, U10, and V10. Table 1 compares the RMSE of ClimaX and the baselines. ClimaX achieves the lowest error for all variables, and the performance gap with the best baseline UNet is larger at increasing lead times. ClimaX also has significant performance gains over its scratch-trained counterpart Cli-ViT, showing the effectiveness of our pretraining procedure in capturing features that are generally useful for various temporal prediction tasks.

### 4.3. Climate projection

To further test the generality of ClimaX, we evaluate the model on ClimateBench [WP+22], a recent benchmark designed for testing machine learning models for climate projections. The goal of ClimateBench is to predict the annual mean global distributions of surface temperature, diurnal temperature range, precipitation, and the 90th percentile of precipitation, given the four anthropogenic forcing factors: carbon dioxide ($CO_2$), sulfur dioxide ($SO_2$), black carbon (BC), and methane ($CH_4$). We note that this is not a temporal modeling task,
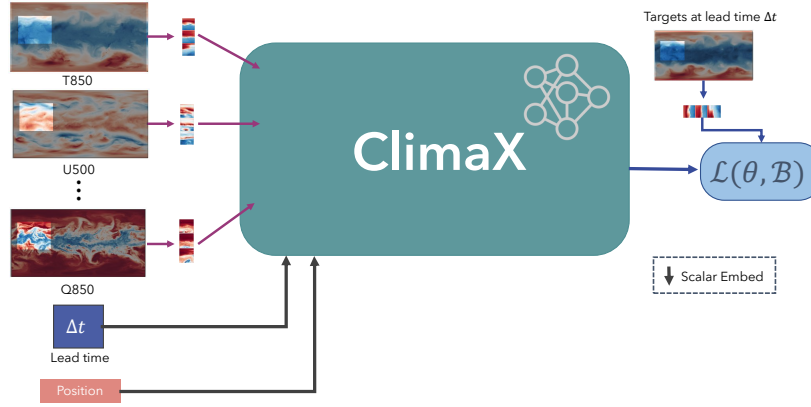
**Figure 7:** Finetuning setup for Regional Forecasting in North America.
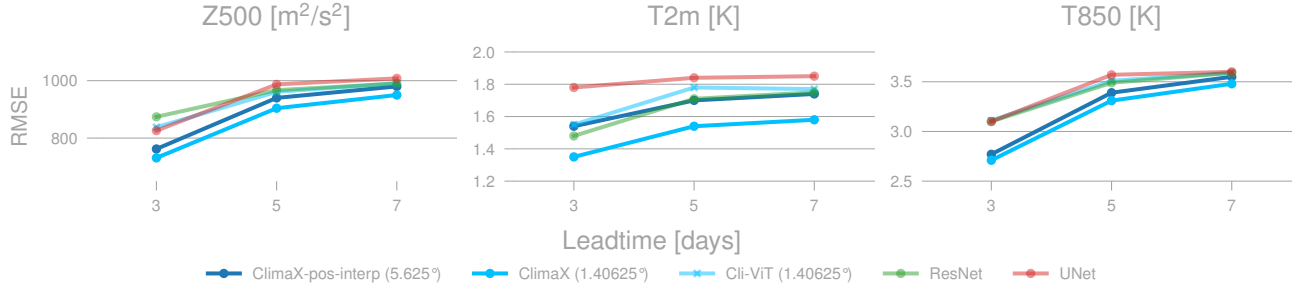


**Figure 8:** Performance on Regional (North America) forecasting for key variables.

as we do not predict the future given the past. Instead, we answer questions like *what will be the annual mean temperature for a specified $CO_2$ level?* In particular, note that the input variables and the task itself are completely different from pretraining.

Figure 9 illustrates the finetuning pipeline of ClimaX for ClimateBench. As the input and output variables are unseen during pretraining, we replace the pretrained embedding layers and prediction heads with newly initialized networks, while keeping the attention layers and the variable aggregation module. We consider two finetuning protocols, in which we either freeze[4] (ClimaX$_{frozen}$) or finetune (ClimaX) the attention layers. In addition, we introduce two components to the pipeline in Figure 2. We use a history of the preceding ten years of the forcing factors to make predictions for a particular year, creating an input of shape $T \times V \times H \times W$. Each time slice of the input goes through variable tokenization, variable aggregation, and the attention layers as usual, which output a feature tensor of shape $T \times h \times w \times D$, where $D$ is the embedding size. The feature tensor then goes through a global average pooling layer, reducing the dimension to $T \times D$. Finally, the 10-year history is aggregated using a cross-attention layer before being fed to the prediction head, which linearly transforms the $D$-dimensional feature vector to a $H \times W$ map. The history aggregation and the global pooling modules are the two additions to the original ClimaX architecture. These architectural designs are inspired by the neural network baseline in [WP+22].

We compare ClimaX with ClimaX$_{frozen}$, Cli-ViT, and the best baseline from ClimateBench. Following [WP+22], we use the standard mean squared error (Equation (1) without the weighting term) as the loss function. We evaluate all methods on RMSE, NRMSE$_s$ (Spatial), NRMSE$_g$ (Global), and Total = NRMSE$_s$ + 5 × NRMSE$_g$ [WP+22]. Details of the metrics are in Appendix D.1. Table 2 shows the results. ClimaX$_{frozen}$ performs the best in predicting two temperature-related variables, followed by ClimaX. This shows that

---

[4]We finetune the LayerNorm in ClimaX$_{frozen}$, as suggested by Lu, Grover, et al. [Lu+22].

**Table 1:** RMSE of ClimaX and baselines on 5.625° ERA5-S2S prediction tasks.

| | T850 | | T2m | | U10 | | V10 | |
|---|---|---|---|---|---|---|---|---|
| | Weeks 3-4 | Weeks 5-6 | Weeks 3-4 | Weeks 5-6 | Weeks 3-4 | Weeks 5-6 | Weeks 3-4 | Weeks 5-6 |
| Resnet | 2.12 | 2.13 | 1.88 | 2.16 | 1.91 | 1.94 | 1.52 | 1.59 |
| Unet | 1.91 | 1.95 | 1.67 | 1.79 | 1.85 | 1.90 | 1.52 | 1.57 |
| Cli-ViT | 1.96 | 1.96 | 1.79 | 1.90 | 1.83 | 1.92 | 1.51 | 1.56 |
| ClimaX | **1.89** | **1.92** | **1.66** | **1.70** | **1.81** | **1.86** | **1.50** | **1.54** |



**Figure 9:** Finetuning pipeline for ClimateBench. A different set of input and output variables requires different embedding layers and prediction heads. Attention layers can be frozen or finetuned.

the pretrained attention layers can serve as a strong feature extractor in seemingly unrelated tasks. Where downstream data is scarce (ClimateBench has only 754 data points), further finetuning the attention layer can lead to overfitting and thus slightly hurt the performance. In two precipitation-related tasks, ClimaX$_{frozen}$ slightly underperforms ClimateBench baseline in terms of NRMSE$_s$ and NRMSE$_g$ but outperforms on RMSE. We hypothesize that this was because ClimaX did not observe the precipitation variable during pretraining, which has very different behaviors from other variables.

## 4.4. Climate model downscaling

Climate models are often run at coarse grids due to their high computational cost. Although these predictions are useful in understanding large-scale climate trends, they do not provide sufficient detail to analyze regional and local phenomena. Downscaling aims to obtain higher-resolution projections and reduce biases from the outputs of these models. To evaluate the applicability of ClimaX to the task of climate model downscaling, we construct a new dataset based on CMIP6 and ERA5 data sources for coarse inputs and higher resolution targets. Specifically, we use all MPI-ESM, a dataset from CMIP6, and its variables listed in Table 8 at 5.625° as input, and train separate models to downscale to each ERA5 target variable at 1.40625°. We compare ClimaX with Cli-ViT and the two CNN baselines, UNet and ResNet, as most recent deep downscaling methods [Van+17; Rod+18; Höh+20; VKG19; LGD20] are based on convolution. We were not able to compare with YNet [LGD20], the current best method on deep downscaling as we did not have access to high-resolution auxiliary data such as elevation and topographical information. For all methods, we first bilinearly interpolate the input to match the resolution of the desired output before feeding it to the model. We evaluate all methods on RMSE, Pearson correlation, and Mean bias, which were commonly used in existing deep downscaling works [Van+17; LGD20]. Details of the metrics are in Appendix D.1.

Table 3 compares ClimaX and the baselines quantitatively. ClimaX achieves the lowest RMSE and a mean

**Table 2:** Performance of ClimaX and the baselines on ClimateBench. Spatial and Global denote the normalized root mean squared error $\text{NRMSE}_s$ and the NRMSE of the global mean $\text{NRMSE}_g$, respectively. Total is a weighted combination of Spatial and Global.

| | Surface temperature | | | | Diurnal temperature range | | | | Precipitation | | | | 90th percentile precipitation | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Spatial | Global | Total | RMSE | Spatial | Global | Total | RMSE | Spatial | Global | Total | RMSE | Spatial | Global | Total | RMSE |
| ClimateBench-NN (reproduced) | 0.123 | 0.080 | 0.524 | 0.404 | 7.465 | 1.233 | 13.632 | 0.150 | 2.349 | **0.151** | 3.104 | 0.553 | 3.108 | **0.282** | 4.517 | 1.594 |
| ClimateBench-NN (paper) | 0.107 | 0.044 | 0.327 | N/A | 9.917 | 1.372 | 16.778 | N/A | **2.128** | 0.209 | 3.175 | N/A | **2.610** | 0.346 | **4.339** | N/A |
| Cli-ViT | 0.086 | 0.044 | 0.305 | 0.362 | 6.997 | 1.759 | 15.792 | 0.146 | 2.224 | 0.241 | 3.430 | 0.550 | 2.800 | 0.329 | 4.447 | 1.579 |
| ClimaX | 0.086 | **0.043** | 0.300 | 0.362 | 7.148 | 0.961 | 11.952 | 0.147 | 2.360 | 0.206 | 3.390 | 0.554 | 2.739 | 0.332 | 4.397 | 1.575 |
| ClimaX$_{\text{frozen}}$ | **0.085** | **0.043** | **0.297** | **0.360** | **6.688** | **0.810** | **10.739** | **0.144** | 2.193 | 0.183 | 3.110 | **0.549** | 2.681 | 0.342 | 4.389 | **1.572** |

**Table 3:** Performance of ClimaX and the baselines on downscaling from MPI-ESM (5.625°) to ERA5 (1.40625°).

| | Z500 | | | T850 | | | T2m | | | U10 | | | V10 | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | RMSE | Pearson | Mean bias | RMSE | Pearson | Mean bias | RMSE | Pearson | Mean bias | RMSE | Pearson | Mean bias | RMSE | Pearson | Mean bias |
| ResNet | 825.75 | 0.96 | −108.54 | 3.60 | 0.96 | 0.19 | 2.89 | 0.98 | 0.14 | 4.05 | 0.65 | 0.06 | 4.11 | 0.45 | 0.09 |
| UNet | 858.35 | 0.95 | 35.10 | 3.66 | 0.96 | −0.34 | 2.95 | 0.98 | 0.16 | 4.09 | 0.64 | −0.06 | 4.13 | 0.44 | 0.08 |
| Cli-ViT | 811.61 | 0.96 | −54.32 | 3.58 | 0.97 | −0.29 | 2.80 | 0.99 | −0.06 | 4.01 | 0.66 | −0.08 | 4.07 | 0.47 | 0.01 |
| ClimaX | **807.43** | 0.96 | **2.70** | **3.49** | 0.97 | **−0.11** | **2.79** | 0.99 | **−0.06** | **3.99** | 0.66 | **0.04** | **4.06** | 0.47 | −0.02 |

bias closest to 0 for all three target variables, and performs similarly to the baselines in terms of Pearson correlation. While pretrained to perform forecasting, ClimaX has successfully captured the spatial structure of weather data, which helps in downstream tasks like downscaling. Figure 10 visualizes the downscaled predictions of ClimaX for the three target variables. The input is at a much lower resolution and contains a lot of bias compared to the ground truth. While the prediction is missing some fine details, it has successfully captured the general structure of the ERA5 data and removed input biases.

## 4.5. Scaling laws analysis

Transformers have shown favorable scaling properties for language [Kap+20; Hof+22], vision [Zha+22a], or even multi-modal tasks [Hen+20b; Hen+21; Ree+22b]. That is, their performance improves with respect to data size and model capacity given sufficient compute. In this section, we study the scaling laws of ClimaX in weather forecasting. Figure 11 presents the performance of ClimaX as a function of data size and model capacity. The $x$-axis is the pretraining data size measured in Gigabytes, which corresponds to 1 to 5 CMIP6 datasets, and the $y$-axis shows the RMSE of ClimaX on the 3-day forecasting task. We compare four ClimaX models with different capacities by varying the embedding dimension from 128 to 1024. All experiments are conducted on the 5.625° data. The error rate of the two biggest models decreases consistently as we increase the data and model size. This highlights the unique ability of ClimaX in learning from diverse and heterogeneous data sources, which allows us to further improve the performance by simply pretraining on more data. However, the two smaller models do not scale as well as the bigger ones, where increasing data size does not gain much improvement or can sometimes hurt performance. This result shows that larger models not only perform better but are also more data efficient.

In addition to data size and model capacity, data resolution is another important scaling dimension in the context of weather and climate. In many vision tasks such as classification, understanding the general, high-level structure of the image is sufficient to make accurate predictions. To model the underlying complex physical processes that govern weather and climate, however, it is important for a model to look at fine-grained details of the input in order to understand the spatial and temporal structure of data as well as the interactions between different variables. High-resolution data contains finer details and local processes of weather conditions that are not present in the low-resolution data, and thus provides stronger signals for training deep learning models. Figure 12 compares the performance of ClimaX pretrained and finetuned on 5.625° and 1.40625° data on global forecasting. Except for T2m at 1 day and 3 days lead times, ClimaX (1.40625°) consistently achieves lower RMSE and higher ACC than the low-resolution model. We note that for the high-resolution data we have to use a larger patch size (4 compared to 2 for low-resolution data) due to lack of memory issue. We can further improve the performance of ClimaX on the 1.40625° data by reducing the patch size, as the
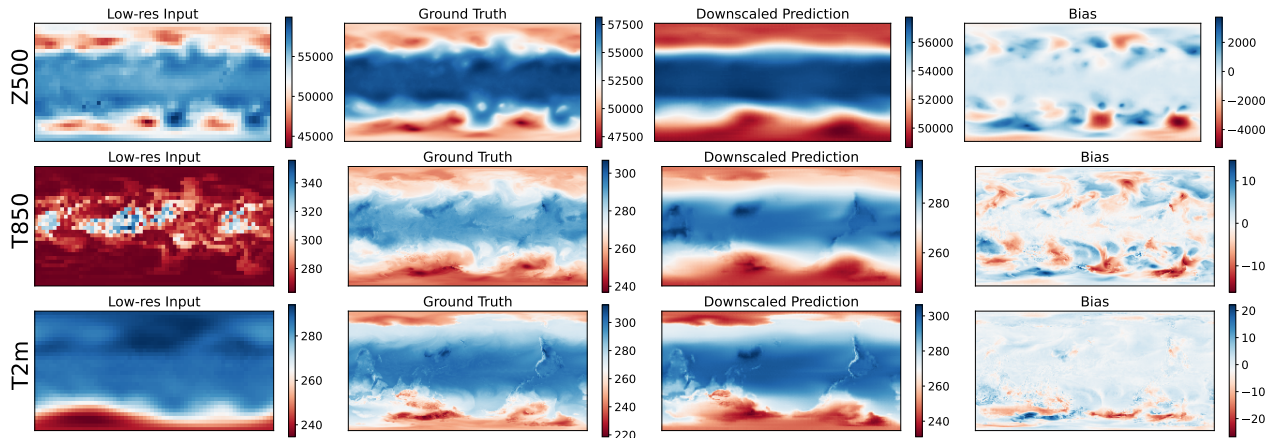
**Figure 10:** Example visualizations of downscaled prediction of key variables by ClimaX.
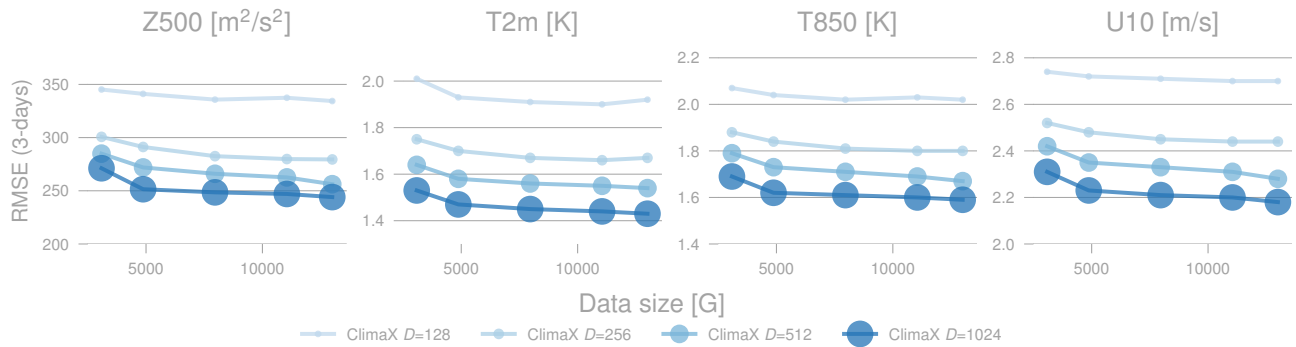


**Figure 11:** Error on ERA5 3-day forecasting for different variables with respect to CMIP6 5.625° data seen during pre-training. Bigger models are more sample efficient.

model is able to capture better details.

### 4.6. Ablation studies

In the main forecasting results, we finetune a separate ClimaX model for each target variable at each lead time, as we found this protocol led to the best performance. However, this can be computationally expensive, as finetuning cost scales linearly with respect to the number of target variables and lead times. In this section, we consider different finetuning alternatives to investigate the trade-off between computation and performance.

### 4.6.1. Should we finetune ClimaX for each variable separately or all at once?

Instead of finetuning ClimaX for each target variable separately, we could alternatively finetune once to predict all variables in the input simultaneously, which we denote as ClimaX-all-vars. Figure 13 shows that ClimaX-all-vars achieves comparable performance to ClimaX in most of the tasks and only underperforms for forecasting T2m. This suggests that with a limited budget, one can finetune ClimaX to predict all target variables at the same time without losing much performance.

### 4.6.2. Should we do iterative forecast or direct forecast?

To avoid finetuning a different model for each lead time, we can finetune ClimaX to make predictions at a short horizon such as 6 hours, and roll out the predictions during inference to make forecasts at longer
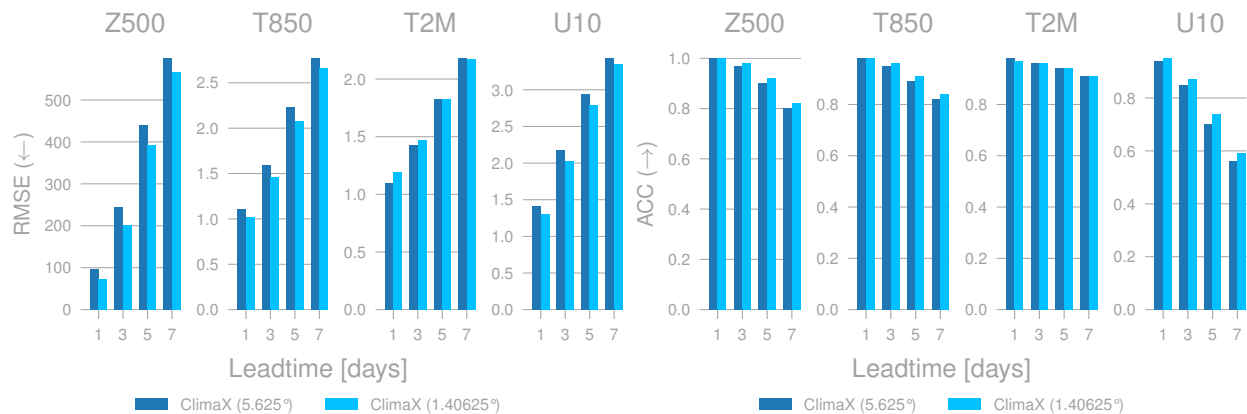
**Figure 12:** Scaling performance with respect to data resolution. Despite a larger patch size, ClimaX (1.40625°) achieves consistently better performance than the low-resolution model on almost all tasks, except for T2m forecast at 1 day and 3 days lead times.

horizons. We call this model ClimaX-iter, where *iter* stands for iterative prediction [Ras+20]. We note that in order to roll out more than one step, ClimaX-iter must predict for all input variables, or in other words. This provides the benefit of finetuning a single model that can predict for any target variable at any lead time. Figure 13 shows that ClimaX-iter works reasonably well up to 1-day prediction, but the performance degrades significantly at longer lead times. This is not surprising, because ClimaX-iter is not finetuned to predict multiple steps into the future, leading to quick error accumulation. One can employ a multi-step objective for finetuning as in Pathak, Subramanian, et al. [Pat+22] to achieve better results.

### 4.6.3. Can we finetune ClimaX to work for all lead times?

Another way to avoid finetuning for each lead time separately is to finetune a lead-time-conditioned model. Specifically, during finetuning, we randomize the lead time from 6 hours to 7 days, resembling the pretraining setting. Note that unlike ClimaX-iter, we still have to finetune a separate model for each target variable. We call this model ClimaX-cont, wherein *cont* stands for *continuous*, a standard term used in previous works [Ras+20]. Figure 13 shows that ClimaX-cont performs competitively on 6-hour to 7-day forecasting, but fails to extrapolate to 2 weeks and 1 month lead times that are unseen during training. One can also randomize the lead time from 6 hours to 1 month, but that means the model sees much fewer data points for each target lead time, potentially hurting the performance.

The cost for finetuning each set of weights is a constant $C$, which is about 15 hours on an $8 \times \text{V100}s$. Among different finetuning protocols, ClimaX is the most expensive, whose total cost is $C \times \#variables \times \#lead\_times$, scaling linearly with the number of target variables and lead times. Following ClimaX are ClimaX-all-vars and ClimaX-cont, whose total costs are $C \times \#lead\_times$ and $C \times \#variables$, respectively. Finally, ClimaX-iter is the cheapest finetuning protocol, where we only have to finetune a single model that works for all target variables and at all lead times. The performance is proportional to the computational cost, as ClimaX is the best performing model, while ClimaX-iter is the worst.

## 5. Discussion and Future Work

The scaling of datasets, model architectures, and computation has resulted in a transformative impact in various subdisciplines of artificial intelligence, from natural language and speech processing to computer vision, as well as scientific applications in biology and chemistry. In particular, it has led to the emergence of general-purpose foundation models that are trained on large datasets and compute clusters, and can be easily adapted to a variety of downstream tasks efficiently, both in terms of compute and data supervision. Our work
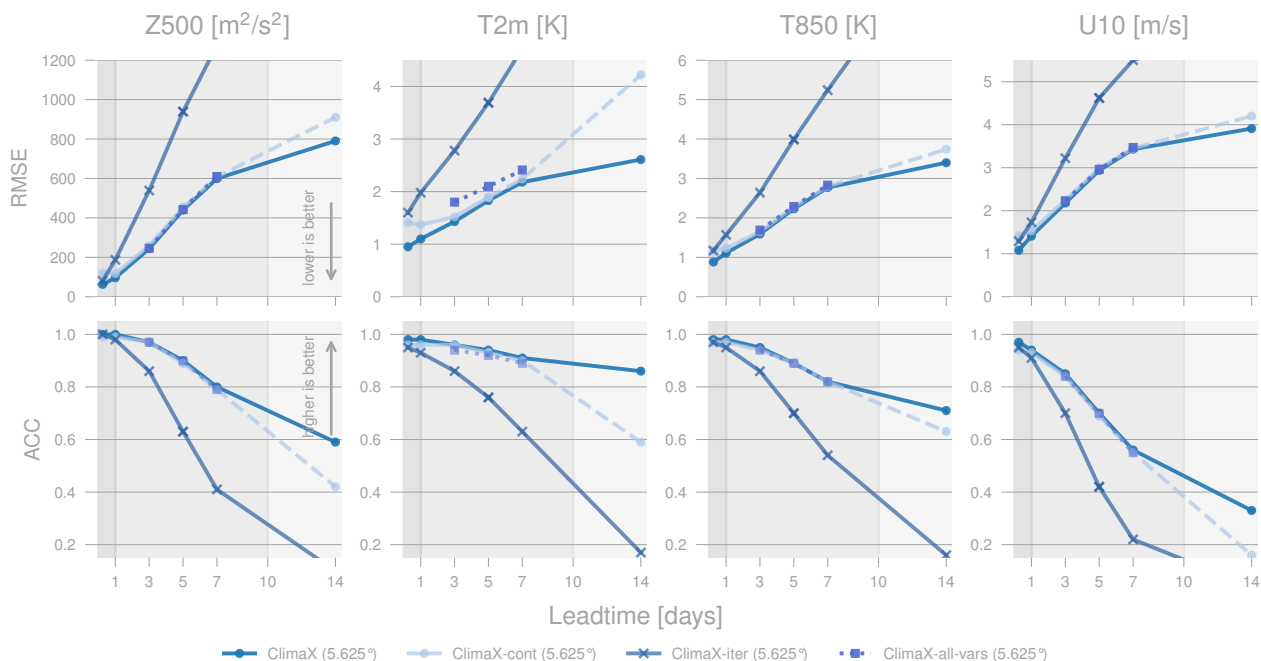
**Figure 13:** Performance of ClimaX and its variations on weather forecasting. ClimaX-cont is a lead-time-conditioned model that we finetune to make predictions at 6 hours to 7 days. ClimaX-iter forecasts at a 6-hour lead time and rolls out the predictions to forecast at longer horizons. ClimaX-all-vars predicts the future conditions of all variables in the input at particular lead-times.

represents a pioneering effort to enable such broad scaling and generality in data-driven models for weather and climate. This approach goes beyond the limitations of both traditional numerical modeling and existing data-driven forecasting methods. Unlike ClimaX, numerical models scale only in terms of computation and not in terms of dataset size, whereas existing data-driven models are typically limited to specific tasks and lack general-purpose applicability across a wide range of tasks.

In addition to traditional considerations in language and vision, foundation models like ClimaX open up new opportunities for scaling through the use of simulation datasets and grid resolutions. To simplify our approach, we chose to use pretraining datasets that include standard variables that have been benchmarked in previous research on data-driven forecasting [Ras+20; Pat+22]. Additionally, we avoided datasets that simulate future scenarios under different forcings to prevent any potential leakage for the climate projection task. Future research could explore incorporating both observational and simulated datasets that include a wider range of climate variables, higher spatiotemporal resolutions, and even extend into future scenarios. Further, we showed that resolution plays a crucial role in scaling of ClimaX. Due to our compute restrictions, we trained ClimaX on low to moderate resolutions. Nevertheless, our empirical trends suggest that scaling to higher resolutions (0.25°) is likely to lead to even better results.

Scaling efforts in the future can benefit from better sequence modeling architectures, especially those designed for multimodal spatiotemporal inputs. As we saw in ClimaX, the number of channels for climate datasets is much greater than those handled for standard multimodal settings (e.g., audio-video, vision-language models). Moreover, in practice, there is also a significant range of resolutions across different climate datasets. This heterogeneity drastically increases the raw length of input sequences for standard architectures such as ViT. In the future, we believe that investigating single multi-scale architectures (e.g., [Fan+21]) can potentially aid in scaling to such diverse multi-resolution and multi-modal datasets by learning to infer features relevant to atmospheric phenomena at increasing spatial resolutions.

In conclusion, we believe that the generality of our approach has potential applications beyond the tasks considered in this work. It would be interesting to explore the generalization of a pretrained ClimaX backbone to other Earth systems science tasks, such as predicting extreme weather events [Mir+19; Sil+17] and assessing anthropogenic contributions to climate change [Ros+08; HT13], as well as broader domains that are closely tied to weather and climate conditions, such as agriculture, demography, and actuarial sciences.

## Acknowledgments

## References

[Ado14]     IPCC Adopted. "Climate change 2014 synthesis report." In: *IPCC: Geneva, Szwitzerland* (2014).

[Arc+20]     Troy Arcomano, Istvan Szunyogh, Jaideep Pathak, Alexander Wikner, Brian R Hunt, and Edward Ott. "A machine learning-based global atmospheric forecast model." In: *Geophysical Research Letters* 47.9 (2020), e2020GL087776.

[Bal+22]     V Balaji, Fleur Couvreux, Julie Deshayes, Jacques Gautrais, Frédéric Hourdin, and Catherine Rio. "Are general circulation models obsolete?" In: *Proceedings of the National Academy of Sciences* 119.47 (2022), e2202075119.

[Bau+20]     Peter Bauer, Tiago Quintino, Nils Wedi, Antonio Bonanni, Marcin Chrust, Willem Deconinck, Michail Diamantakis, Peter Düben, Stephen English, Johannes Flemming, et al. *The ecmwf scalability programme: Progress and plans.* European Centre for Medium Range Weather Forecasts, 2020.

[BGS20]     Lea Beusch, Lukas Gudmundsson, and Sonia I Seneviratne. "Emulating Earth system model temperatures with MESMER: from global mean temperature trajectories to grid-point-level realizations on land." In: *Earth System Dynamics* 11.1 (2020), pp. 139–159.

[Bi+22]     Kaifeng Bi, Lingxi Xie, Hengheng Zhang, Xin Chen, Xiaotao Gu, and Qi Tian. "Pangu-Weather: A 3D High-Resolution Model for Fast and Accurate Global Weather Forecast." In: *arXiv preprint arXiv:2211.02556* (2022).

[BMMG20]     Jorge Baño-Medina, Rodrigo Manzanas, and José Manuel Gutiérrez. "Configuration and inter-comparison of deep learning neural models for statistical downscaling." In: *Geoscientific Model Development* 13.4 (2020), pp. 2109–2124.

[Bom+21]     Rishi Bommasani, Drew A. Hudson, et al. "On the Opportunities and Risks of Foundation Models." In: *ArXiv* (2021). URL: https://crfm.stanford.edu/assets/report.pdf.

[Bou+10]     Philippe Bougeault, Zoltan Toth, Craig Bishop, Barbara Brown, David Burridge, De Hui Chen, Beth Ebert, Manuel Fuentes, Thomas M Hamill, Ken Mylne, et al. "The THORPEX interactive grand global ensemble." In: *Bulletin of the American Meteorological Society* 91.8 (2010), pp. 1059–1072.

[Bra+22]     Johannes Brandstetter, Rianne van den Berg, Max Welling, and Jayesh K Gupta. "Clifford Neural Layers for PDE Modeling." In: *arXiv preprint arXiv:2209.04934* (2022).

[Bro+20]     Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. "Language models are few-shot learners." In: *Advances in neural information processing systems* 33 (2020), pp. 1877–1901.

[BTB15]     Peter Bauer, Alan Thorpe, and Gilbert Brunet. "The quiet revolution of numerical weather prediction." In: *Nature* 525.7567 (2015), pp. 47–55.

[BWW22]     Johannes Brandstetter, Daniel Worrall, and Max Welling. "Message Passing Neural PDE Solvers." In: *arXiv preprint arXiv:2202.03376* (2022).

[Cho+22]     Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. "PaLM: Scaling language modeling with pathways." In: *arXiv preprint arXiv:2204.02311* (2022).

[Con+22]     Yezhen Cong, Samar Khanna, Chenlin Meng, Patrick Liu, Erik Rozi, Yutong He, Marshall Burke, David B Lobell, and Stefano Ermon. "Satmae: Pre-training transformers for temporal and multi-spectral satellite imagery." In: *arXiv preprint arXiv:2207.08051* (2022).

[DB18]     Peter D Dueben and Peter Bauer. "Challenges and design choices for global weather and climate models based on machine learning." In: *Geoscientific Model Development* 11.10 (2018), pp. 3999–4009.

[Dev+18]     Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. "Bert: Pre-training of deep bidirectional transformers for language understanding." In: *arXiv preprint arXiv:1810.04805* (2018).

[Dos+20]     Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. "An image is worth 16x16 words: Transformers for image recognition at scale." In: *arXiv preprint arXiv:2010.11929* (2020).

[Ern21]     Lukas Ernst. "Structured Attention Transformers on Weather Prediction." MA thesis. ETH Zurich, Scalable Parallel Computing Laboratory, 2021.

[Eyr+16]     Veronika Eyring, Sandrine Bony, Gerald A Meehl, Catherine A Senior, Bjorn Stevens, Ronald J Stouffer, and Karl E Taylor. "Overview of the Coupled Model Intercomparison Project Phase 6 (CMIP6) experimental design and organization." In: *Geoscientific Model Development* 9.5 (2016), pp. 1937–1958.

[Fan+21]     Haoqi Fan, Bo Xiong, Karttikeya Mangalam, Yanghao Li, Zhicheng Yan, Jitendra Malik, and Christoph Feichtenhofer. "Multiscale vision transformers." In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2021, pp. 6824–6835.

[GB22]     Jayesh K Gupta and Johannes Brandstetter. "Towards Multi-spatiotemporal-scale Generalized PDE Modeling." In: *arXiv preprint arXiv:2209.15616* (2022).

[GKH15]     Aditya Grover, Ashish Kapoor, and Eric Horvitz. "A deep hybrid model for weather forecasting." In: *Proceedings of the 21th ACM SIGKDD international conference on knowledge discovery and data mining*. 2015, pp. 379–386.

[Gro22]     Aditya Grover. "Rethinking Machine Learning for Climate Science: A Dataset Perspective." In: *AAAI Symposium on The Role of AI in Responding to Climate Challenges*. 2022.

[Har+20]     Charles R. Harris, K. Jarrod Millman, Stéfan J. van der Walt, Ralf Gommers, Pauli Virtanen, David Cournapeau, Eric Wieser, Julian Taylor, Sebastian Berg, Nathaniel J. Smith, Robert Kern, Matti Picus, Stephan Hoyer, Marten H. van Kerkwijk, Matthew Brett, Allan Haldane, Jaime Fernández del Río, Mark Wiebe, Pearu Peterson, Pierre Gérard-Marchant, Kevin Sheppard, Tyler Reddy, Warren Weckesser, Hameer Abbasi, Christoph Gohlke, and Travis E. Oliphant. "Array programming with NumPy." In: *Nature* 585.7825 (Sept. 2020), pp. 357–362. DOI: 10.1038/s41586-020-2649-2. URL: https://doi.org/10.1038/s41586-020-2649-2.

[He+16]     Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. "Deep residual learning for image recognition." In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2016, pp. 770–778.

[He+22]      Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. "Masked autoencoders are scalable vision learners." In: *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 2022, pp. 16000–16009.

[Hen+20a]    Dan Hendrycks, Xiaoyuan Liu, Eric Wallace, Adam Dziedzic, Rishabh Krishnan, and Dawn Song. "Pretrained Transformers Improve Out-of-Distribution Robustness." In: *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. 2020, pp. 2744–2751.

[Hen+20b]    Tom Henighan, Jared Kaplan, Mor Katz, Mark Chen, Christopher Hesse, Jacob Jackson, Heewoo Jun, Tom B Brown, Prafulla Dhariwal, Scott Gray, et al. "Scaling laws for autoregressive generative modeling." In: *arXiv preprint arXiv:2010.14701* (2020).

[Hen+21]     Lisa Anne Hendricks, John Mellor, Rosalia Schneider, Jean-Baptiste Alayrac, and Aida Nematzadeh. "Decoupling the Role of Data, Attention, and Losses in Multimodal Transformers." In: *Transactions of the Association for Computational Linguistics* 9 (2021), pp. 570–585.

[Her+18]     H Hersbach, B Bell, P Berrisford, G Biavati, A Horányi, J Muñoz Sabater, J Nicolas, C Peubey, R Radu, I Rozum, et al. "ERA5 hourly data on single levels from 1979 to present." In: *Copernicus Climate Change Service (C3S) Climate Data Store (CDS)* 10 (2018).

[Her+20]     Hans Hersbach, Bill Bell, Paul Berrisford, Shoji Hirahara, András Horányi, Joaquín Muñoz-Sabater, Julien Nicolas, Carole Peubey, Raluca Radu, Dinand Schepers, et al. "The ERA5 global reanalysis." In: *Quarterly Journal of the Royal Meteorological Society* 146.730 (2020), pp. 1999–2049.

[HH17]       Stephan Hoyer and Joe Hamman. "xarray: N-D labeled Arrays and Datasets in Python." In: *Journal of Open Research Software* 5.1 (Apr. 2017), p. 10. DOI: 10.5334/jors.148.

[Hof+22]     Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al. "Training Compute-Optimal Large Language Models." In: *arXiv preprint arXiv:2203.15556* (2022).

[Höh+20]     Kevin Höhlein, Michael Kern, Timothy Hewson, and Rüdiger Westermann. "A comparative study of convolutional neural network models for wind field downscaling." In: *Meteorological Applications* 27.6 (2020), e1961.

[HT13]       Mikael Höök and Xu Tang. "Depletion of fossil fuels and anthropogenic climate change—A review." In: *Energy policy* 52 (2013), pp. 797–809.

[Hua+16]     Gao Huang, Yu Sun, Zhuang Liu, Daniel Sedra, and Kilian Q Weinberger. "Deep networks with stochastic depth." In: *European conference on computer vision*. Springer. 2016, pp. 646–661.

[Hun+19]     Chris Huntingford, Elizabeth S Jeffers, Michael B Bonsall, Hannah M Christensen, Thomas Lees, and Hui Yang. "Machine learning and artificial intelligence to aid climate change research and preparedness." In: *Environmental Research Letters* 14.12 (2019), p. 124007.

[Hur+13]     James W Hurrell, Marika M Holland, Peter R Gent, Steven Ghan, Jennifer E Kay, Paul J Kushner, J-F Lamarque, William G Large, D Lawrence, Keith Lindsay, et al. "The community earth system model: a framework for collaborative research." In: *Bulletin of the American Meteorological Society* 94.9 (2013), pp. 1339–1360.

[Hwa+19]     Jessica Hwang, Paulo Orenstein, Judah Cohen, Karl Pfeiffer, and Lester Mackey. "Improving subseasonal forecasting in the western US with machine learning." In: *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 2019, pp. 2325–2335.

[Kal03]      Eugenia Kalnay. *Atmospheric modeling, data assimilation and predictability*. Cambridge university press, 2003.

[Kap+20]     Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. "Scaling laws for neural language models." In: *arXiv preprint arXiv:2001.08361* (2020).

[Kas+21]     K Kashinath, M Mustafa, A Albert, JL Wu, C Jiang, S Esmaeilzadeh, K Azizzadenesheli, R Wang, A Chattopadhyay, A Singh, et al. "Physics-informed machine learning: case studies for weather and climate modelling." In: *Philosophical Transactions of the Royal Society A* 379.2194 (2021), p. 20200093.

[KB14]       Diederik P Kingma and Jimmy Ba. "Adam: A method for stochastic optimization." In: *arXiv preprint arXiv:1412.6980* (2014).

[Kei22]      Ryan Keisler. "Forecasting Global Weather with Graph Neural Networks." In: *arXiv preprint arXiv:2202.07575* (2022).

[Koc+21]     Dmitrii Kochkov, Jamie A Smith, Ayya Alieva, Qing Wang, Michael P Brenner, and Stephan Hoyer. "Machine learning–accelerated computational fluid dynamics." In: *Proceedings of the National Academy of Sciences* 118.21 (2021), e2101784118.

[Lam+22]     Remi Lam, Alvaro Sanchez-Gonzalez, Matthew Willson, Peter Wirnsberger, Meire Fortunato, Alexander Pritzel, Suman Ravuri, Timo Ewalds, Ferran Alet, Zach Eaton-Rosen, et al. "GraphCast: Learning skillful medium-range global weather forecasting." In: *arXiv preprint arXiv:2212.12794* (2022).

[LGD20]      Yumin Liu, Auroop R Ganguly, and Jennifer Dy. "Climate downscaling using YNet: A deep convolutional network with skip connections and fusion." In: *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining.* 2020, pp. 3145–3153.

[LH17]       Ilya Loshchilov and Frank Hutter. "Decoupled weight decay regularization." In: *arXiv preprint arXiv:1711.05101* (2017).

[Li+20]      Zongyi Li, Nikola Kovachki, Kamyar Azizzadenesheli, Burigede Liu, Kaushik Bhattacharya, Andrew Stuart, and Anima Anandkumar. "Fourier neural operator for parametric partial differential equations." In: *arXiv preprint arXiv:2010.08895* (2020).

[Liu+21]     Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. "Swin transformer: Hierarchical vision transformer using shifted windows." In: *Proceedings of the IEEE/CVF International Conference on Computer Vision.* 2021, pp. 10012–10022.

[Liu+22]     Ze Liu, Han Hu, Yutong Lin, Zhuliang Yao, Zhenda Xie, Yixuan Wei, Jia Ning, Yue Cao, Zheng Zhang, Li Dong, Furu Wei, and Baining Guo. "Swin Transformer V2: Scaling Up Capacity and Resolution." In: *International Conference on Computer Vision and Pattern Recognition (CVPR).* 2022.

[Lor67]      Edward Lorenz. "The nature and theory of the general circulation of the atmosphere." In: *World meteorological organization* 161 (1967).

[LSZ15]      Kody Law, Andrew Stuart, and Konstantinos Zygalakis. "Data assimilation." In: *Cham, Switzerland: Springer* 214 (2015), p. 52.

[Lu+21]      Lu Lu, Pengzhan Jin, Guofei Pang, Zhongqiang Zhang, and George Em Karniadakis. "Learning nonlinear operators via DeepONet based on the universal approximation theorem of operators." In: *Nature Machine Intelligence* 3.3 (2021), pp. 218–229.

[Lu+22]      Kevin Lu, Aditya Grover, Pieter Abbeel, and Igor Mordatch. "Pretrained transformers as universal computation engines." In: *AAAI Conference on Artificial Intelligence.* 2022.

[Lyn08]      Peter Lynch. "The origins of computer weather prediction and climate modeling." In: *Journal of computational physics* 227.7 (2008), pp. 3431–3444.

[Man+20]     Laura A Mansfield, Peer J Nowack, Matt Kasoar, Richard G Everitt, William J Collins, and Apostolos Voulgarakis. "Predicting global patterns of long-term climate change from short-term simulations using machine learning." In: *npj Climate and Atmospheric Science* 3.1 (2020), pp. 1–9.

[Mar+22]     Stratis Markou, James Requeima, Wessel P Bruinsma, Anna Vaughan, and Richard E Turner. "Practical Conditional Neural Processes Via Tractable Dependent Predictions." In: *arXiv preprint arXiv:2203.08775* (2022).

[MD+21]      Valérie Masson-Delmotte, Panmao Zhai, Anna Pirani, Sarah L Connors, Clotilde Péan, Sophie Berger, Nada Caud, Y Chen, L Goldfarb, MI Gomis, et al. "Climate change 2021: the physical science basis." In: *Contribution of working group I to the sixth assessment report of the intergovernmental panel on climate change* 2 (2021).

[Mee+00]     Gerald A Meehl, George J Boer, Curt Covey, Mojib Latif, and Ronald J Stouffer. "The coupled model intercomparison project (CMIP)." In: *Bulletin of the American Meteorological Society* 81.2 (2000), pp. 313–318.

[Mir+19]     Diego G Miralles, Pierre Gentine, Sonia I Seneviratne, and Adriaan J Teuling. "Land–atmospheric feedbacks during droughts and heatwaves: state of the science and current challenges." In: *Annals of the New York Academy of Sciences* 1436.1 (2019), pp. 19–35.

[Mou+23a]    Soukayna Mouatadid, Paulo Orenstein, Genevieve Flaspohler, Judah Cohen, Miruna Oprescu, Ernest Fraenkel, and Lester Mackey. "Adaptive bias correction for improved subseasonal forecasting." In: *Nature Communications* 14.1 (2023), p. 3482.

[Mou+23b]    Soukayna Mouatadid, Paulo Orenstein, Genevieve Elaine Flaspohler, Miruna Oprescu, Judah Cohen, Franklyn Wang, Sean Edward Knight, Maria Geogdzhayeva, Samuel James Levang, Ernest Fraenkel, et al. "SubseasonalClimateUSA: A Dataset for Subseasonal Forecasting and Benchmarking." In: *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*. 2023.

[Ore+19]     Boris N Oreshkin, Dmitri Carpov, Nicolas Chapados, and Yoshua Bengio. "N-BEATS: Neural basis expansion analysis for interpretable time series forecasting." In: *arXiv preprint arXiv:1905.10437* (2019).

[Pas+19a]    Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. "PyTorch: An Imperative Style, High-Performance Deep Learning Library." In: *Advances in Neural Information Processing Systems (NeurIPS)*. Curran Associates, Inc., 2019, pp. 8024–8035.

[Pas+19b]    Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. "Pytorch: An imperative style, high-performance deep learning library." In: *Advances in neural information processing systems* 32 (2019).

[Pat+22]     Jaideep Pathak, Shashank Subramanian, Peter Harrington, Sanjeev Raja, Ashesh Chattopadhyay, Morteza Mardani, Thorsten Kurth, David Hall, Zongyi Li, Kamyar Azizzadenesheli, et al. "Fourcastnet: A global data-driven high-resolution weather model using adaptive fourier neural operators." In: *arXiv preprint arXiv:2202.11214* (2022).

[Phi56]      Norman A Phillips. "The general circulation of the atmosphere: A numerical experiment." In: *Quarterly Journal of the Royal Meteorological Society* 82.352 (1956), pp. 123–164.

[Pro+18]     Liudmila Prokhorenkova, Gleb Gusev, Aleksandr Vorobev, Anna Veronika Dorogush, and Andrey Gulin. "CatBoost: unbiased boosting with categorical features." In: *Advances in neural information processing systems* 31 (2018).

[Rad+21]     Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. "Learning transferable visual models from natural language supervision." In: *International Conference on Machine Learning*. PMLR. 2021, pp. 8748–8763.

[Ram+22]   Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. "Hierarchical text-conditional image generation with clip latents." In: *arXiv preprint arXiv:2204.06125* (2022).

[Ras+20]   Stephan Rasp, Peter D Dueben, Sebastian Scher, Jonathan A Weyn, Soukayna Mouatadid, and Nils Thuerey. "WeatherBench: a benchmark data set for data-driven weather forecasting." In: *Journal of Advances in Modeling Earth Systems* 12.11 (2020), e2020MS002203.

[Rav+21]   Suman Ravuri, Karel Lenc, Matthew Willson, Dmitry Kangin, Remi Lam, Piotr Mirowski, Megan Fitzsimons, Maria Athanassiadou, Sheleem Kashem, Sam Madge, et al. "Skilful precipitation nowcasting using deep generative models of radar." In: *Nature* 597.7878 (2021), pp. 672–677.

[Ree+22a]   Colorado J Reed, Ritwik Gupta, Shufan Li, Sarah Brockman, Christopher Funk, Brian Clipp, Salvatore Candido, Matt Uyttendaele, and Trevor Darrell. "Scale-MAE: A Scale-Aware Masked Autoencoder for Multiscale Geospatial Representation Learning." In: *arXiv preprint arXiv:2212.14532* (2022).

[Ree+22b]   Scott Reed, Konrad Zolna, Emilio Parisotto, Sergio Gómez Colmenarejo, Alexander Novikov, Gabriel Barth-maron, Mai Giménez, Yury Sulsky, Jackie Kay, Jost Tobias Springenberg, Tom Eccles, Jake Bruce, Ali Razavi, Ashley Edwards, Nicolas Heess, Yutian Chen, Raia Hadsell, Oriol Vinyals, Mahyar Bordbar, and Nando de Freitas. "A Generalist Agent." In: *Transactions on Machine Learning Research* (2022). Featured Certification. URL: https://openreview.net/forum?id=1ikK0kHjvj.

[Rei+19]   Markus Reichstein, Gustau Camps-Valls, Bjorn Stevens, Martin Jung, Joachim Denzler, Nuno Carvalhais, et al. "Deep learning and process understanding for data-driven Earth system science." In: *Nature* 566.7743 (2019), pp. 195–204.

[RFB15]   Olaf Ronneberger, Philipp Fischer, and Thomas Brox. "U-Net: Convolutional networks for biomedical image segmentation." In: *International Conference on Medical image computing and computer-assisted intervention*. Springer. 2015, pp. 234–241.

[Rod+18]   Eduardo Rocha Rodrigues, Igor Oliveira, Renato Cunha, and Marco Netto. "DeepDownscale: a deep learning strategy for high-resolution weather forecast." In: *2018 IEEE 14th International Conference on e-Science (e-Science)*. IEEE. 2018, pp. 415–422.

[Roh22]   Ankit Rohatgi. *Webplotdigitizer: Version 4.6*. 2022. URL: https://automeris.io/WebPlotDigitizer.

[Ros+08]   Cynthia Rosenzweig, David Karoly, Marta Vicarelli, Peter Neofotis, Qigang Wu, Gino Casassa, Annette Menzel, Terry L Root, Nicole Estrella, Bernard Seguin, et al. "Attributing physical and biological impacts to anthropogenic climate change." In: *Nature* 453.7193 (2008), pp. 353–357.

[RRH22]   AR Ravishankara, David A Randall, and James W Hurrell. "Complex and yet predictable: The message of the 2021 Nobel Prize in Physics." In: *Proceedings of the National Academy of Sciences* 119.2 (2022), e2120669119.

[RT21]   Stephan Rasp and Nils Thuerey. "Data-driven medium-range weather prediction with a resnet pretrained on climate simulations: A new model for weatherbench." In: *Journal of Advances in Modeling Earth Systems* 13.2 (2021), e2020MS002405.

[Sac+18]   DA Sachindra, Khandakar Ahmed, Md Mamunur Rashid, S Shahid, and BJC Perera. "Statistical downscaling of precipitation using machine learning techniques." In: *Atmospheric research* 212 (2018), pp. 240–258.

[Sat04]   Masaki Satoh. *Atmospheric circulation dynamics and circulation models*. Springer Science & Business Media, 2004.

[Sch+17]   Tapio Schneider, Shiwei Lan, Andrew Stuart, and Joao Teixeira. "Earth system modeling 2.0: A blueprint for models that learn from observations and targeted high-resolution simulations." In: *Geophysical Research Letters* 44.24 (2017), pp. 12–396.

[Sch18]   Sebastian Scher. "Toward data-driven weather and climate forecasting: Approximating a simple general circulation model with deep learning." In: *Geophysical Research Letters* 45.22 (2018), pp. 12–616.

[Sch+21]   Martin G Schultz, Clara Betancourt, Bing Gong, Felix Kleinert, Michael Langguth, Lukas Hubert Leufen, Amirpasha Mozaffari, and Scarlet Stadtler. "Can deep learning beat numerical weather prediction?" In: *Philosophical Transactions of the Royal Society A* 379.2194 (2021), p. 20200097.

[Sil+17]   Jana Sillmann, Thordis Thorarinsdottir, Noel Keenlyside, Nathalie Schaller, Lisa V Alexander, Gabriele Hegerl, Sonia I Seneviratne, Robert Vautard, Xuebin Zhang, and Francis W Zwiers. "Understanding, modeling and predicting weather and climate extremes: Challenges and opportunities." In: *Weather and climate extremes* 18 (2017), pp. 65–74.

[SM19]   Sebastian Scher and Gabriele Messori. "Weather and climate forecasting with neural networks: using general circulation models (GCMs) with different complexity as a study ground." In: *Geoscientific Model Development* 12.7 (2019), pp. 2797–2809.

[Søn+20]   Casper Kaae Sønderby, Lasse Espeholt, Jonathan Heek, Mostafa Dehghani, Avital Oliver, Tim Salimans, Shreya Agrawal, Jason Hickey, and Nal Kalchbrenner. "MetNet: A neural weather model for precipitation forecasting." In: *arXiv preprint arXiv:2003.12140* (2020).

[Tao+20]   Rohan Taori, Achal Dave, Vaishaal Shankar, Nicholas Carlini, Benjamin Recht, and Ludwig Schmidt. "Measuring robustness to natural distribution shifts in image classification." In: *Advances in Neural Information Processing Systems* 33 (2020), pp. 18583–18599.

[TL18]   Sean J Taylor and Benjamin Letham. "Forecasting at scale." In: *The American Statistician* 72.1 (2018), pp. 37–45.

[Tou+21]   Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. "Training data-efficient image transformers & distillation through attention." In: *International Conference on Machine Learning*. PMLR. 2021, pp. 10347–10357.

[TP12]   Sebastian Thrun and Lorien Pratt. *Learning to learn.* Springer Science & Business Media, 2012.

[Van+17]   Thomas Vandal, Evan Kodra, Sangram Ganguly, Andrew Michaelis, Ramakrishna Nemani, and Auroop R Ganguly. "Deepsd: Generating high resolution climate change projections through single image super-resolution." In: *Proceedings of the 23rd acm sigkdd international conference on knowledge discovery and data mining*. 2017, pp. 1663–1672.

[Vas+17]   Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. "Attention is all you need." In: *Advances in neural information processing systems* 30 (2017).

[Vau+21]   Anna Vaughan, Will Tebbutt, J Scott Hosking, and Richard E Turner. "Convolutional conditional neural processes for local climate downscaling." In: *arXiv preprint arXiv:2101.07950* (2021).

[Ver+22]   Robert Verkuil, Ori Kabeli, Yilun Du, Basile IM Wicky, Lukas F Milles, Justas Dauparas, David Baker, Sergey Ovchinnikov, Tom Sercu, and Alexander Rives. "Language models generalize beyond natural proteins." In: *bioRxiv* (2022), pp. 2022–12.

[Vit+22]   F. Vitart, A. W. Robertson, A. Spring, F. Pinault, R. Roškar, W. Cao, S. Bech, A. Bienkowski, N. Caltabiano, E. De Coning, B. Denis, A. Dirkson, J. Dramsch, P. Dueben, J. Gierschendorf, H. S. Kim, K. Nowak, D. Landry, L. Lledó, L. Palma, S. Rasp, and S. Zhou. "Outcomes of the WMO Prize Challenge to Improve Subseasonal to Seasonal Predictions Using Artificial Intelligence." In: *Bulletin of the American Meteorological Society* 103.12 (Dec. 2022), E2878–E2886. DOI: 10.1175/bams-d-22-0046.1.

[VKG19]   Thomas Vandal, Evan Kodra, and Auroop R Ganguly. "Intercomparison of machine learning methods for statistical downscaling: the case of daily and extreme precipitation." In: *Theoretical and Applied Climatology* 137.1 (2019), pp. 557–570.

[VR18]   Frédéric Vitart and Andrew W Robertson. "The sub-seasonal to seasonal prediction project (S2S) and the prediction of extreme events." In: *npj Climate and Atmospheric Science* 1.1 (2018), pp. 1–7.

[Wan+22]  Wenhui Wang, Hangbo Bao, Li Dong, Johan Bjorck, Zhiliang Peng, Qiang Liu, Kriti Aggarwal, Owais Khan Mohammed, Saksham Singhal, Subhojit Som, et al. "Image as a foreign language: Beit pretraining for all vision and vision-language tasks." In: *arXiv preprint arXiv:2208.10442* (2022).

[War10]  Thomas Tomkins Warner. *Numerical weather and climate prediction*. cambridge university press, 2010.

[WDC20]  Jonathan A Weyn, Dale R Durran, and Rich Caruana. "Improving data-driven global weather prediction using deep convolutional neural networks on a cubed sphere." In: *Journal of Advances in Modeling Earth Systems* 12.9 (2020), e2020MS002109.

[Web+20]  Theodore Weber, Austin Corotan, Brian Hutchinson, Ben Kravitz, and Robert Link. "Deep learning for creating surrogate models of precipitation in Earth system models." In: *Atmospheric Chemistry and Physics* 20.4 (2020), pp. 2303–2317.

[Wed+15]  NP Wedi, P Bauer, W Denoninck, M Diamantakis, M Hamrud, C Kuhnlein, S Malardel, K Mogensen, G Mozdzynski, and PK Smolarkiewicz. *The modelling infrastructure of the Integrated Forecasting System: Recent advances and future challenges*. European Centre for Medium-Range Weather Forecasts, 2015.

[Wey+21]  Jonathan A Weyn, Dale R Durran, Rich Caruana, and Nathaniel Cresswell-Clay. "Sub-seasonal forecasting with a large ensemble of deep-learning weather prediction models." In: *Journal of Advances in Modeling Earth Systems* 13.7 (2021), e2021MS002502.

[Wig19]  Ross Wightman. *PyTorch Image Models*. https://github.com/rwightman/pytorch-image-models. 2019. DOI: 10.5281/zenodo.4414861.

[WP+22]  Duncan Watson-Parris, Yuhan Rao, Dirk Olivié, Øyvind Seland, Peer Nowack, Gustau Camps-Valls, Philip Stier, Shahine Bouabid, Maura Dewey, Emilie Fons, et al. "ClimateBench v1. 0: A Benchmark for Data-Driven Climate Projections." In: *Journal of Advances in Modeling Earth Systems* 14.10 (2022), e2021MS002954.

[WW97]  Robert L Wilby and Thomas ML Wigley. "Downscaling general circulation model output: a review of methods and limitations." In: *Progress in physical geography* 21.4 (1997), pp. 530–548.

[YL20]  Yuan Yuan and Lei Lin. "Self-supervised pretraining of transformers for satellite image time series classification." In: *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing* 14 (2020), pp. 474–487.

[Yua+21]  Lu Yuan, Dongdong Chen, Yi-Ling Chen, Noel Codella, Xiyang Dai, Jianfeng Gao, Houdong Hu, Xuedong Huang, Boxin Li, Chunyuan Li, et al. "Florence: A new foundation model for computer vision." In: *arXiv preprint arXiv:2111.11432* (2021).

[Zha+19]  Fuqing Zhang, Y Qiang Sun, Linus Magnusson, Roberto Buizza, Shian-Jiann Lin, Jan-Huey Chen, and Kerry Emanuel. "What is the predictability limit of midlatitude weather?" In: *Journal of the Atmospheric Sciences* 76.4 (2019), pp. 1077–1091.

[Zha+22a]  Xiaohua Zhai, Alexander Kolesnikov, Neil Houlsby, and Lucas Beyer. "Scaling vision transformers." In: *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 2022, pp. 12104–12113.

[Zha+22b]  Chongzhi Zhang, Mingyuan Zhang, Shanghang Zhang, Daisheng Jin, Qiang Zhou, Zhongang Cai, Haiyu Zhao, Xianglong Liu, and Ziwei Liu. "Delving deep into the generalization of vision transformers under distribution shifts." In: *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 2022, pp. 7277–7286.

[Zho+21]  Haoyi Zhou, Shanghang Zhang, Jieqi Peng, Shuai Zhang, Jianxin Li, Hui Xiong, and Wancai Zhang. "Informer: Beyond efficient transformer for long sequence time-series forecasting." In: *Proceedings of the AAAI conference on artificial intelligence*. Vol. 35. 12. 2021, pp. 11106–11115.

[Zhu18]  J Zhuang. *xESMF: Universal regridder for geospatial data*. 2018.

# A. Model

This section presents the implementation details and hyperparameters of ClimaX and the two CNN baselines UNet and ResNet.

## A.1. ClimaX

### A.1.1. Implementation details

ClimaX receives a tensor of shape $V \times H \times W$ and outputs a tensor of shape $V' \times H \times W$, where the number of input and output variables $V$ and $V'$ can vary between different datasets[5]. To do that, we assume a set $\mathcal{V}$ that contains all possible variables we could encounter during pretraining and finetuning. Each variable in $\mathcal{V}$ has a separate token embedding layer.

The variable tokenization module tokenizes the input to a sequence of $V \times h \times w$ tokens, with each token being a vector of size $p^2$. After that, for each token, we extract the corresponding embedding layer that transforms the token to a vector of dimension $D$. Each embedding layer is a single convolution layer with $in\_channels = 1, out\_channels = D, kernel\_size = p, stride = p$. This results in a tensor of shape $V \times h \times w \times D$.

To differentiate between tokens of different input variables, we add the sequence with a *variable positional embedding*, which is a tensor of shape $|\mathcal{V}| \times D$. For each input variable, we extract the corresponding variable positional embedding to add to its tokens. After that, all tokens go through the variable aggregation module, which outputs a tensor of shape $h \times w \times D$.

The tokens are then fed to the attention layers, which output a tensor of the same shape $h \times w \times D$. The prediction head takes each token of dimension $D$ and maps it to a vector of dimension $|\mathcal{V}| \times p^2$, and the output is reshaped to $|\mathcal{V}| \times H \times W$. Finally, we extract predictions of $V'$ target variables and compute the loss.

### A.1.2. Hyperparameters

**Table 4:** Default hyperparameters of ClimaX

| Hyperparameter | Meaning | Value |
|---|---|---|
| $\mathcal{V}$ | Default variables | All ERA5 variables in Table 9 |
| $|\mathcal{V}|$ | Number of default variables | 48 |
| $p$ | Patch size | 2 for 5.625° <br> 4 for 1.40625° |
| $D$ | Embedding dimension | 1024 |
| Depth | Number of ViT blocks | 8 |
| # heads | Number of attention heads | 16 |
| 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 | 1024 |
| Drop path | For stochastic depth [Hua+16] | 0.1 |
| Dropout | Dropout rate | 0.1 |

---

[5]The spatial resolution $H \times W$ can also vary. In that case, we employ the common practice of interpolating the positional embedding, and everything else remains the same [Dos+20; Tou+21].

### A.2. CNN Baselines

### A.2.1. ResNet Hyperparameters

We use the following hyperparameters for ResNet in all of our experiments.

**Table 5:** 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 |

### A.2.2. UNet Hyperparameters

We borrow our UNet implementation from PDEArena [GB22]. We use the following hyperparameters for UNet in all of our experiments.

**Table 6:** 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 |
| Channel multiplications | Determine the number of output channels for Down and Up blocks | $[1, 2, 2, 4]$ |
| Blocks | Number of blocks | 2 |
| Use attention | If use attention in Down and Up blocks | False |
| Dropout | Dropout rate | 0.1 |

### A.2.3. Other implementation details

Following the implementation of ResNet in Rasp, Dueben, et al. [Ras+20], Rasp and Thuerey [RT21], and Ernst [Ern21], we found the following details important for the performance of both CNN baselines:

- Use Batch normalization
- Use Leakyrelu with a slope of 0.3 as the activation function
- Postnorm instead of Prenorm
- Use periodic convolutions in the longitude direction but not the latitude direction.
- Use a kernel size of 7 in the first CNN layer.

## B. Training details

**Data normalization** We normalized all inputs during pre-training as well as fine-tuning. For each variable, at each pressure level (for atmospheric variables), we compute the mean and standard deviation to normalize them to zero mean and unit variance. We de-normalize the predictions to get back to the original range before computing evaluation metrics.

**Software and hardware stack** We use PyTorch [Pas+19a], `timm` [Wig19], `numpy` [Har+20] and `xarray` [HH17] to manage our data and model training. We used 32GB NVIDIA V100 devices for training. For pretraining we distribute the batch across 80 V100s on AzureML. We leverage `fp16` floating point precision in our model.

### B.1. Pretraining

### B.1.1. Objective

We use the loss function in Equation (1) for pretraining.

### B.1.2. Optimization

We used the AdamW optimizer [KB14; LH17] with parameters $(\beta_1 = 0.9, \beta_2 = 0.95)$. We used weight decay of $1e - 5$ for all parameters except for the positional embedding. We used a learning rate of $5e - 4$, with a linear warmup schedule for 10000 steps (5 epochs), followed by a cosine-annealing schedule for 190000 steps (95 epochs).

### B.2. Finetuning

### B.2.1. Objective

We use lat-weighted MSE in Equation (1) for finetuning ClimaX in temporal forecasting and downscaling tasks. In ClimateBench, we finetune using standard MSE without the weighting term, as this led to better results and was suggested by [WP+22].

### B.2.2. Optimization

For all tasks, we used AdamW with parameters $(\beta_1 = 0.9, \beta_2 = 0.999)$. We used weight decay of $1e - 5$ for all parameters except for the positional embedding. We used a linear warmup schedule for 10000 steps (5 epochs), followed by a cosine-annealing schedule for 90000 steps (45 epochs). The learning rate for each task is as follows:

**Table 7:** Learning rate for finetuning ClimaX in different downstream tasks

| Task | Learning rate |
|---|---|
| Weather forecasting | $5e - 7$ |
| Climate projection | $5e - 4$ |
| Climate downscaling | $5e - 5$ |

We used a small learning rate for weather forecasting as the task resembles pretraining. For downscaling, we used a larger learning rate, as the nature of the task is different from pretraining, even though the input variables are similar. In climate projection, we needed to initialize new weights for the embedding layers and prediction heads, and thus used a similar learning rate to training from scratch.

## C. Datasets

### C.1. CMIP6-ClimaX

We created CMIP6-ClimaX for pretraining ClimaX, which consists of 5 datasets from the CMIP6 project. We downloaded the datasets from the official CMIP6 search interface at https://esgf-data.dkrz.de/search/cmip6-dkrz/. These datasets share the following attributes:

- Experiment ID: historical
- Table ID: 6hrPlevPt, i.e., 6-hourly data on pressure levels.
- Variant label: r1i1p1f1. The variant label distinguishes among closely related simulations by a single model, in which "r" specifies the initial condition, "i" specifies the observational dataset and initialization method used for determining the initial condition, "p" specifies the perturbed physics version of the model, and "f" specifies the forcing index.

All datasets have a temporal coverage from 1850 to 2015 and a temporal resolution of 6 hours. We chose these datasets as they contain similar climate variables at similar vertical levels to ERA5. We note that there are more than 5 datasets from CMIP6 that suit our selection criteria, but we were not able to download others due to some issues on the data servers. We regridded these datasets to 5.625° and 1.40625° using the xesmf Python package [Zhu18] using bilinear interpolation. We provide a detailed description of these 5 data sources and the available variables we used to construct CMIP6-ClimaX in Table 8.

We note that AWI and HAMMOZ are not the best data sources for higher resolution 1.40625° training, because their original resolution at 250 km is lower than 1.40625°, which is about 156 km. We wanted to use other higher-resolution datasets but were not able to download them. We believe pretraining on other high-resolution datasets would lead to better performance.

## C.2. ERA5

We use the preprocessed version of ERA5 from WeatherBench [Ras+20] for finetuning ClimaX. WeatherBench was created as a standard benchmark data and evaluation framework for comparing data-driven weather forecasting models. WeatherBench regridded the original ERA5 at 0.25° to three lower resolutions: 5.625°, 2.8125°, and 1.40625°. See https://confluence.ecmwf.int/display/CKB/ERA5%3A+data+documentation for more details of the raw ERA5 data. Table 9 summarizes the variables we use for finetuning ClimaX.

### C.2.1. ERA5-NA

We constructed ERA5-NA from ERA5 to evaluate ClimaX and the baselines on regional forecasting. ERA-NA has the same set of variables as in Table 9, but only contains data that belongs to the North America region. To do this, we first identified the latitude and longitude range to form a rectangular area that encapsulates North America, using the standard CORDEX domains https://cordex.org/wp-content/uploads/2012/11/CORDEX-domain-description_231015.pdf. For each data sample, we then extracted the spatial positions that fall into this range, forming in ERA5-NA.

### C.2.2. ERA-S2S

We built ERA5-S2S from ERA5 to serve as a benchmark dataset for sub-seasonal to seasonal prediction. ERA5-S2S consists of two sub-datasets, whose the goals are to predict the biweekly average statistics of target variables in weeks 3 and 4, and weeks 5 and 6, respectively. The input includes all variables in Table 9, while the output variables are are averaged over two weeks, starting from the start of week 3 (5) and to the end of week 4 (6).

## C.3. ClimateBench

We refer to Watson-Parris, Rao, et al. [WP+22] for complete details of ClimateBench.

**Table 8:** Resolution and variables of CMIP6-ClimaX dataset used for pretraining. *Static* represents variables don't depend on time, *Single* represents surface variables, and *Atmospheric* represents time-varying atmospheric properties at the chosen altitudes.

| Data Source | Original resolution | Variables | | |
|---|---|---|---|---|
| | | Type | Abbrev. | Levels |
| MPI | 100km | Single | t2m | |
| | | Single | u10 | |
| | | Single | v10 | |
| | | Atmospheric | z | 50, 250, 500, 600, 700, 850, 925 |
| | | Atmospheric | u | 50, 250, 500, 600, 700, 850, 925 |
| | | Atmospheric | v | 50, 250, 500, 600, 700, 850, 925 |
| | | Atmospheric | t | 50, 250, 500, 600, 700, 850, 925 |
| | | Atmospheric | q | 50, 250, 500, 600, 700, 850, 925 |
| Tai | 100km | Single | t2m | |
| | | Atmospheric | z | 250, 500, 600, 700, 850, 925 |
| | | Atmospheric | u | 250, 500, 850 |
| | | Atmospheric | v | 250, 500, 850 |
| | | Atmospheric | t | 250, 500, 850 |
| | | Atmopheric | q | 250, 500, 600, 700, 850, 925 |
| AWI | 250km | Single | t2m | |
| | | Single | u10 | |
| | | Single | v10 | |
| | | Atmospheric | z | 50, 250, 500, 600, 700, 850, 925 |
| | | Atmospheric | u | 50, 250, 500, 600, 700, 850, 925 |
| | | Atmospheric | v | 50, 250, 500, 600, 700, 850, 925 |
| | | Atmospheric | t | 50, 250, 500, 600, 700, 850, 925 |
| | | Atmospheric | q | 50, 250, 500, 600, 700, 850, 925 |
| HAMMOZ | 250km | Single | t2m | |
| | | Single | u10 | |
| | | Single | v10 | |
| | | Atmospheric | z | 50, 250, 500, 600, 700, 850, 925 |
| | | Atmospheric | u | 50, 250, 500, 600, 700, 850, 925 |
| | | Atmospheric | v | 50, 250, 500, 600, 700, 850, 925 |
| | | Atmospheric | t | 50, 250, 500, 600, 700, 850, 925 |
| | | Atmospheric | q | 50, 250, 500, 600, 700, 850, 925 |
| CMCC | 100km | Atmospheric | z | 50, 250, 500, 600, 700, 850, 925 |
| | | Atmospheric | u | 50, 250, 500, 600, 700, 850, 925 |
| | | Atmospheric | v | 50, 250, 500, 600, 700, 850, 925 |
| | | Atmospheric | t | 250, 500, 850 |

## D. Quantitative evaluation

### D.1. Metrics

This section presents all evaluation metrics we use in Section 4. For all metrics, we denote $\tilde{X}$ and $X$ as the prediction and ground truth, which have a shape of $N \times H \times W$, where $N$ is the number of forecasts, or the number of test samples, $H \times W$ is the spatial resolution. $L(i)$ is the latitude weighting term to account for

**Table 9:** ECMWF variables used in our ERA5 dataset. *Static* represents variables don't depend on time, *Single* represents surface variables, and *Atmospheric* represents time-varying atmospheric properties at the chosen altitudes.

| Type | Variable name | Abbrev. | ECMWF ID | Levels |
|---|---|---|---|---|
| Static | Land-sea mask | LSM | 172 | |
| Static | Orography | | | |
| Single | 2 metre temperature | T2m | 167 | |
| Single | 10 metre U wind component | U10 | 165 | |
| Single | 10 metre V wind component | V10 | 166 | |
| Atmospheric | Geopotential | Z | 129 | 50, 250, 500, 600, 700, 850, 925 |
| Atmospheric | U wind component | U | 131 | 50, 250, 500, 600, 700, 850, 925 |
| Atmospheric | V wind component | V | 132 | 50, 250, 500, 600, 700, 850, 925 |
| Atmospheric | Temperature | T | 130 | 50, 250, 500, 600, 700, 850, 925 |
| Atmospheric | Specific humidity | Q | 133 | 50, 250, 500, 600, 700, 850, 925 |
| Atmospheric | Relative humidity | R | 157 | 50, 250, 500, 600, 700, 850, 925 |

the non-uniformity in areas of the grid cells. We have removed the time notation for simplicity.

### D.1.1. Weather forecasting metrics

**Root mean square error (RMSE)**

$$\text{RMSE} = \frac{1}{N}\sum_{k=1}^{N}\sqrt{\frac{1}{H\times W}\sum_{i=1}^{H}\sum_{j=1}^{W}L(i)(\tilde{X}_{k,i,j}-X_{k,i,j})^2}. \tag{3}$$

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

$$\text{ACC} = \frac{\sum_{k,i,j}L(i)\tilde{X}'_{k,i,j}X'_{k,i,j}}{\sqrt{\sum_{k,i,j}L(i)\tilde{X}'^2_{k,i,j}\sum_{k,i,j}L(i)X'^2_{k,i,j}}}, \tag{4}$$

$$\tilde{X}' = \tilde{X}' - C, X' = X' - C, \tag{5}$$

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

### D.1.2. Climate projection metrics

**Normalized spatial root mean square error ($\text{NRMSE}_s$)** Normalized spatial root mean square error ($\text{NRMSE}_s$) measures the spatial discrepancy between the temporal mean of the prediction and the temporal mean of the ground truth:

$$\text{NRMSE}_s = \sqrt{\left\langle \left(\frac{1}{N}\sum_{k=1}^{N}\tilde{X} - \frac{1}{N}\sum_{k=1}^{N}X\right)^2\right\rangle} \Big/ \frac{1}{N}\sum_{k=1}^{N}\langle X\rangle, \tag{6}$$

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

$$\langle A\rangle = \frac{1}{H\times W}\sum_{i=1}^{H}\sum_{j=1}^{W}L(i)A_{i,j} \tag{7}$$

**Normalized global root mean square error (NRMSE$_g$)** Normalized global root mean square error (NRMSE$_g$) measures the discrepancy between the global mean of the prediction and the global mean of the ground truth:

$$\text{NRMSE}_g = \sqrt{\frac{1}{N}\sum_{k=1}^{N}\left(\langle\tilde{X}\rangle - \langle X\rangle\right)^2}\bigg/ \frac{1}{N}\sum_{k=1}^{N}\langle X\rangle. \tag{8}$$

**Total normalized root mean square error (TRMSE)** Total normalized root mean square error (TRMSE) is the weighted sum of NRMSE$_s$ and NRMSE$_g$:

$$\text{TRMSE} = \text{NRMSE}_s + \alpha \cdot \text{NRMSE}_g, \tag{9}$$

where $\alpha$ is chosen to be 5 as suggested by Watson-Parris, Rao, et al. [WP+22].

### D.1.3. Climate downscaling metrics

**Root mean square error (RMSE)** This is the same as Equation (3).

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

$$\text{Mean bias} = \frac{1}{N \times H \times W}\sum_{k=1}^{N}\sum_{i=1}^{H}\sum_{j=1}^{W}\tilde{X} - \frac{1}{N \times H \times W}\sum_{k=1}^{N}\sum_{i=1}^{H}\sum_{j=1}^{W}X \tag{10}$$

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

$$\rho_{\tilde{X},X} = \frac{\text{cov}(\tilde{X},X)}{\sigma_{\tilde{X}}\sigma_X} \tag{11}$$

### D.2. Results summary

Table 10 and 11 summarize the global forecasting results of ClimaX and the baselines for all target variables and at all lead times. In addition to IFS and the two CNN-based baselines in the main text, we include FourCastNet [Pat+22], PanguWeather [Bi+22], and GraphCast [Lam+22] for comprehensiveness. We want to emphasize that the results obtained by these methods are not comparable with ClimaX, as they were trained on ERA5 at 0.25°, a much higher resolution compared to 5.625° and 1.40625° data used to train ClimaX. In Section 4.5, we had a discussion on how the performance of ClimaX scales favorably with respect to data resolution. We hope this summary will provide future works with an easier comparison with existing baselines.

In spite of being trained on much lower resolutions, ClimaX outperforms FourCastNet in forecasting Z500, T850, and U10 at lead times from 3 days and beyond, in terms of both RMSE and ACC. For T2m, ClimaX achieves better results at horizons longer than 3 days. PanguWeather performs better than ClimaX on most of the tasks, but the gap between the two methods shrinks and becomes negligible as the lead time increases. ClimaX even outperforms PanguWeather in predicting U10 at 7 days lead times. This is because ClimaX is finetuned to perform direct prediction, which mitigates error accumulation for long horizon prediction. GraphCast achieves the lowest RMSE among all methods, but performs worse in terms of ACC compared to ClimaX and PanguWeather.

**Table 10:** RMSE on global forecasting for different target variables at different lead times. Lower is better.

| Variable | Lead time [hr.] | ClimaX 5.625° | ClimaX 1.40625° | FCN[a] 0.25° | PW[b] 0.25° | GC[c] 0.25° | HRES 0.1 | IFS 5.625° | IFS 1.40625° | ResNet 5.625° | UNet 5.625° |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Z500 $[m^2/s^2]$ | 6 | 62.73 | 49.67 | 37.52 | 15.40 | 16.46 | 24.66 | 26.93 | 26.96 | 47.00 | 53.66 |
| | 24 | 96.19 | 72.76 | 81.31 | 42.23 | 38.77 | 45.90 | 51.01 | 50.96 | 86.60 | 132.65 |
| | 72 | 244.08 | 201.00 | 251.96 | 133.12 | 125.78 | 146.37 | 152.15 | 152.20 | 305.22 | 458.84 |
| | 120 | 440.40 | 392.00 | 483.44 | 295.63 | 271.65 | 316.79 | 331.45 | 331.38 | 614.20 | 721.83 |
| | 168 | 599.43 | 566.00 | 680.00 | 504.90 | 466.53 | 535.93 | 549.01 | 548.96 | 806.59 | 819.39 |
| | 336 | 790.26 | 788.43 | nan | nan | nan | nan | 1011.72 | 1011.56 | 835.55 | 866.40 |
| | 720 | 815.25 | 817.52 | nan | nan | nan | nan | nan | nan | 858.98 | 880.34 |
| T2m [K] | 6 | 0.95 | 1.11 | 0.72 | 0.59 | 0.50 | 0.35 | 0.97 | 0.97 | 0.76 | 0.77 |
| | 24 | 1.10 | 1.19 | 0.95 | 0.72 | 0.62 | 0.66 | 1.02 | 1.02 | 0.91 | 1.11 |
| | 72 | 1.43 | 1.47 | 1.38 | 1.05 | 0.94 | 1.06 | 1.30 | 1.30 | 1.70 | 1.91 |
| | 120 | 1.83 | 1.83 | 1.99 | 1.53 | 1.36 | 1.52 | 1.72 | 1.71 | 2.22 | 2.49 |
| | 168 | 2.18 | 2.17 | 2.54 | 2.06 | 1.88 | 2.06 | 2.24 | 2.23 | 2.66 | 2.66 |
| | 336 | 2.61 | 2.67 | nan | nan | nan | nan | 3.31 | 3.30 | 2.86 | 2.79 |
| | 720 | 2.67 | 2.74 | nan | nan | nan | nan | nan | nan | 2.86 | 2.81 |
| T850 [K] | 6 | 0.88 | 0.84 | 0.52 | 0.42 | 0.28 | 0.33 | 0.69 | 0.69 | 0.70 | 0.80 |
| | 24 | 1.11 | 1.02 | 0.81 | 0.72 | 0.58 | 0.70 | 0.87 | 0.87 | 1.26 | 1.25 |
| | 72 | 1.59 | 1.46 | 1.55 | 1.13 | 1.02 | 1.27 | 1.34 | 1.34 | 1.90 | 2.39 |
| | 120 | 2.23 | 2.08 | 2.47 | 1.78 | 1.63 | 1.96 | 2.01 | 2.01 | 2.86 | 3.23 |
| | 168 | 2.77 | 2.66 | 3.30 | 2.60 | 2.41 | 2.78 | 2.82 | 2.82 | 3.51 | 3.50 |
| | 336 | 3.40 | 3.41 | nan | nan | nan | nan | 4.43 | 4.43 | 3.65 | 3.65 |
| | 720 | 3.47 | 3.49 | nan | nan | nan | nan | nan | nan | 3.69 | 3.73 |
| U10 [m/s] | 6 | 1.08 | 1.04 | 0.55 | 0.46 | 0.37 | 0.58 | 0.80 | 0.79 | 0.86 | 1.02 |
| | 24 | 1.41 | 1.31 | 0.99 | 0.90 | 0.80 | 1.15 | 1.11 | 1.11 | 1.27 | 1.68 |
| | 72 | 2.18 | 2.02 | 2.24 | 1.60 | 1.47 | 1.98 | 1.92 | 1.92 | 2.78 | 3.17 |
| | 120 | 2.94 | 2.79 | 3.41 | 2.52 | 2.36 | 2.95 | 2.89 | 2.89 | 3.63 | 3.93 |
| | 168 | 3.43 | 3.35 | 4.18 | 3.46 | 3.25 | 3.87 | 3.81 | 3.81 | 4.15 | 4.08 |
| | 336 | 3.91 | 3.92 | nan | nan | nan | nan | 5.24 | 5.23 | 4.23 | 4.16 |
| | 720 | 3.96 | 3.97 | nan | nan | nan | nan | nan | nan | 4.29 | 4.22 |

[a] FourCastNet [Pat+22]
[b] PanguWeather [Bi+22]
[c] GraphCast [Lam+22]

**Table 11:** ACC on global forecasting for different target variables at different lead times. Higher is better.

| VARIABLE | LEAD TIME | ClimaX | | FCN[a] | PW[b] | GC[c] | HRES | IFS | | ResNet | UNet |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | [hr.] | 5.625° | 1.40625° | 0.25° | 0.25° | 0.25° | 0.1 | 5.625° | 1.40625° | 5.625° | 5.625° |
| Z500 | 6 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 |
| | 24 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 0.99 |
| | 72 | 0.97 | 0.98 | 0.97 | 0.99 | 0.99 | 0.98 | 0.99 | 0.99 | 0.95 | 0.89 |
| | 120 | 0.90 | 0.92 | 0.89 | 0.96 | 0.94 | 0.92 | 0.95 | 0.95 | 0.79 | 0.69 |
| | 168 | 0.80 | 0.82 | 0.76 | 0.87 | 0.83 | 0.78 | 0.87 | 0.87 | 0.57 | 0.57 |
| | 336 | 0.59 | 0.59 | nan | nan | nan | nan | 0.55 | 0.55 | 0.53 | 0.51 |
| | 720 | 0.55 | 0.55 | nan | nan | nan | nan | nan | nan | 0.49 | 0.49 |
| T2m | 6 | 0.98 | 0.98 | 0.99 | 0.99 | 0.98 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 |
| | 24 | 0.98 | 0.97 | 0.98 | 0.99 | 0.98 | 0.98 | 0.99 | 0.99 | 0.98 | 0.98 |
| | 72 | 0.96 | 0.96 | 0.96 | 0.98 | 0.95 | 0.94 | 0.98 | 0.98 | 0.94 | 0.93 |
| | 120 | 0.94 | 0.94 | 0.92 | 0.95 | 0.90 | 0.88 | 0.96 | 0.96 | 0.90 | 0.88 |
| | 168 | 0.91 | 0.91 | 0.87 | 0.92 | 0.81 | 0.77 | 0.93 | 0.93 | 0.86 | 0.86 |
| | 336 | 0.86 | 0.85 | nan | nan | nan | nan | 0.85 | 0.85 | 0.83 | 0.84 |
| | 720 | 0.85 | 0.84 | nan | nan | nan | nan | nan | nan | 0.83 | 0.83 |
| T850 | 6 | 0.98 | 0.99 | 0.99 | 1.00 | 1.00 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 |
| | 24 | 0.98 | 0.98 | 0.98 | 0.99 | 0.99 | 0.98 | 0.99 | 0.99 | 0.97 | 0.97 |
| | 72 | 0.95 | 0.96 | 0.95 | 0.98 | 0.96 | 0.93 | 0.97 | 0.97 | 0.92 | 0.88 |
| | 120 | 0.89 | 0.91 | 0.87 | 0.94 | 0.89 | 0.84 | 0.93 | 0.94 | 0.82 | 0.75 |
| | 168 | 0.82 | 0.84 | 0.77 | 0.87 | 0.75 | 0.68 | 0.87 | 0.87 | 0.68 | 0.69 |
| | 336 | 0.71 | 0.71 | nan | nan | nan | nan | 0.68 | 0.69 | 0.66 | 0.66 |
| | 720 | 0.69 | 0.68 | nan | nan | nan | nan | nan | nan | 0.64 | 0.64 |
| U10 | 6 | 0.97 | 0.97 | 0.99 | 0.99 | 0.99 | 0.99 | 0.98 | 0.98 | 0.98 | 0.97 |
| | 24 | 0.94 | 0.95 | 0.97 | 0.97 | 0.98 | 0.96 | 0.97 | 0.97 | 0.95 | 0.91 |
| | 72 | 0.85 | 0.87 | 0.85 | 0.92 | 0.93 | 0.88 | 0.89 | 0.89 | 0.74 | 0.65 |
| | 120 | 0.70 | 0.74 | 0.64 | 0.80 | 0.82 | 0.74 | 0.76 | 0.76 | 0.52 | 0.37 |
| | 168 | 0.56 | 0.59 | 0.45 | 0.63 | 0.64 | 0.55 | 0.58 | 0.58 | 0.28 | 0.28 |
| | 336 | 0.33 | 0.32 | nan | nan | nan | nan | 0.21 | 0.21 | 0.19 | 0.22 |
| | 720 | 0.29 | 0.28 | nan | nan | nan | nan | nan | nan | 0.17 | 0.21 |

[a] FourCastNet [Pat+22]
[b] PanguWeather [Bi+22]
[c] GraphCast [Lam+22]

# E. Qualitative evaluation

We qualitatively evaluate the performance of CliMax on global forecasting tasks for all target variables and at all lead times. In each figure, the first column is the initial condition of the target variable, which serves as the input, the second column is the ground truth of the target variable at a particular lead time, the third column is the prediction of ClimaX, and the last column is the bias, which is the difference between the prediction and the ground truth.

## E.1. Nowcasting



**Figure 14:** Example forecasts from ClimaX at 6-hour lead time compared to ground truth ERA5.

## E.2. Short and medium-range weather forecasting



**Figure 15:** Example forecasts from ClimaX at 1-day lead time compared to ground truth ERA5.



**Figure 16:** Example forecasts from ClimaX at 3-day lead time compared to ground truth ERA5.
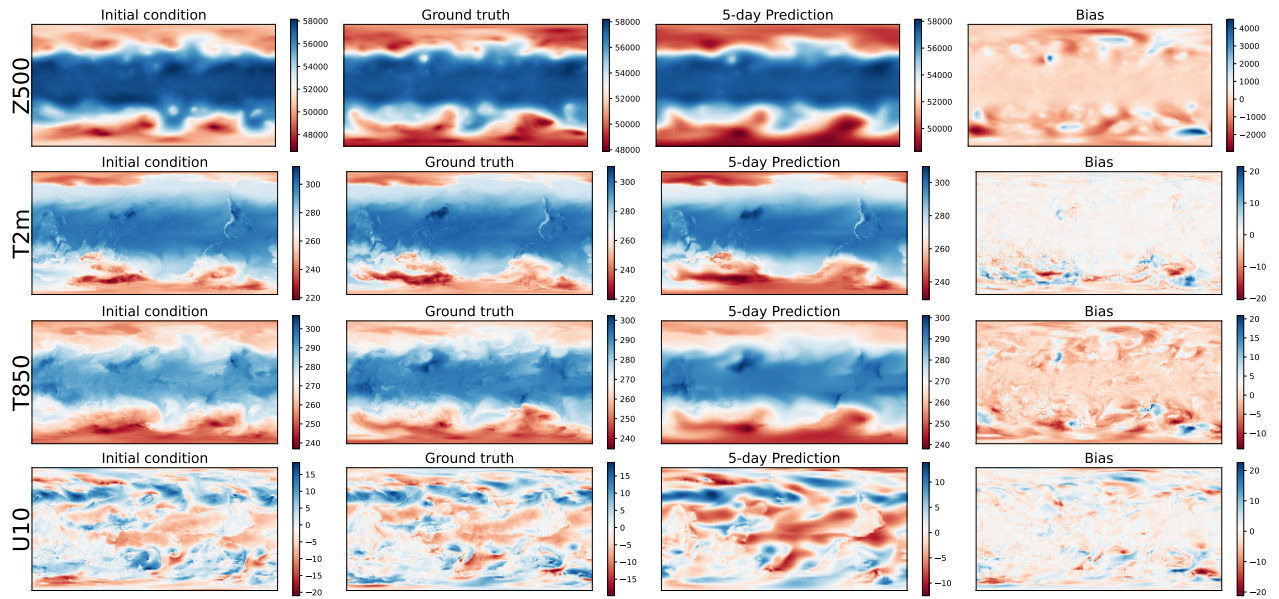
**Figure 17:** Example forecasts from ClimaX at 5-day lead time compared to ground truth ERA5.
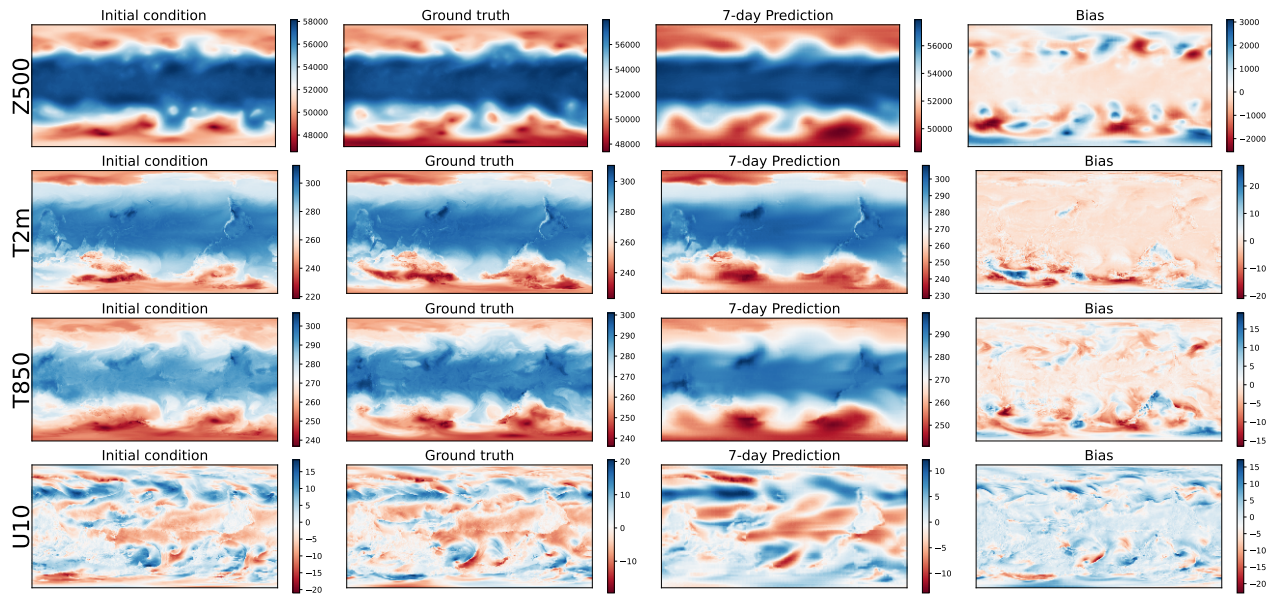


**Figure 18:** Example forecasts from ClimaX at 7-day lead time compared to ground truth ERA5.
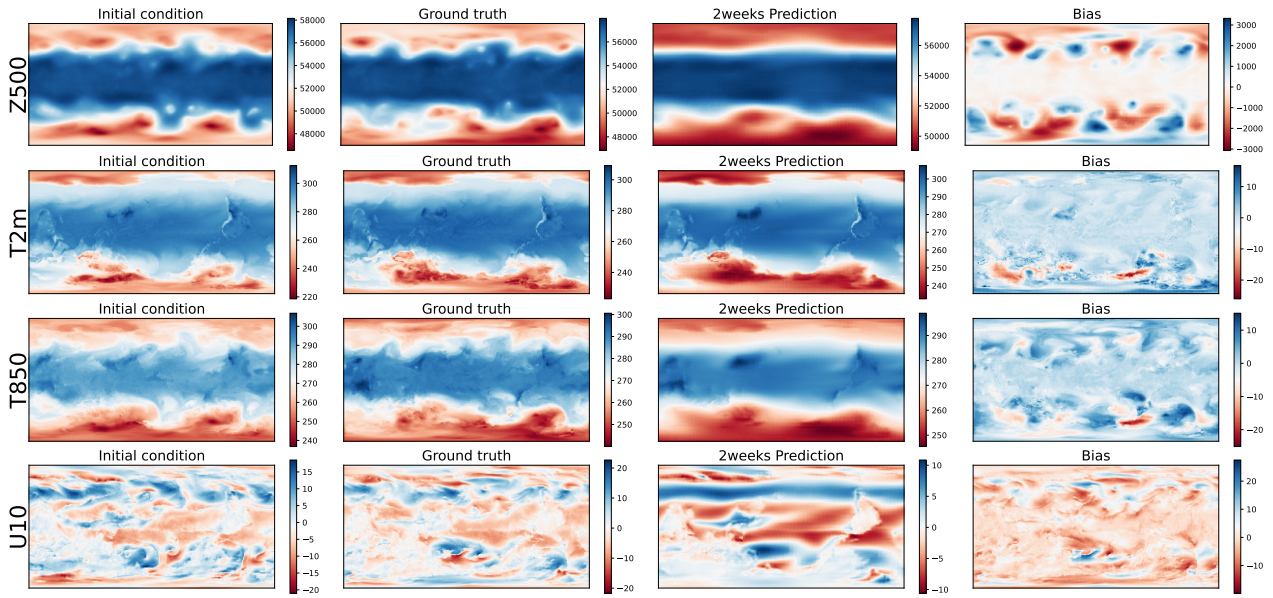
## E.3. Longer horizon instantaneous forecasting



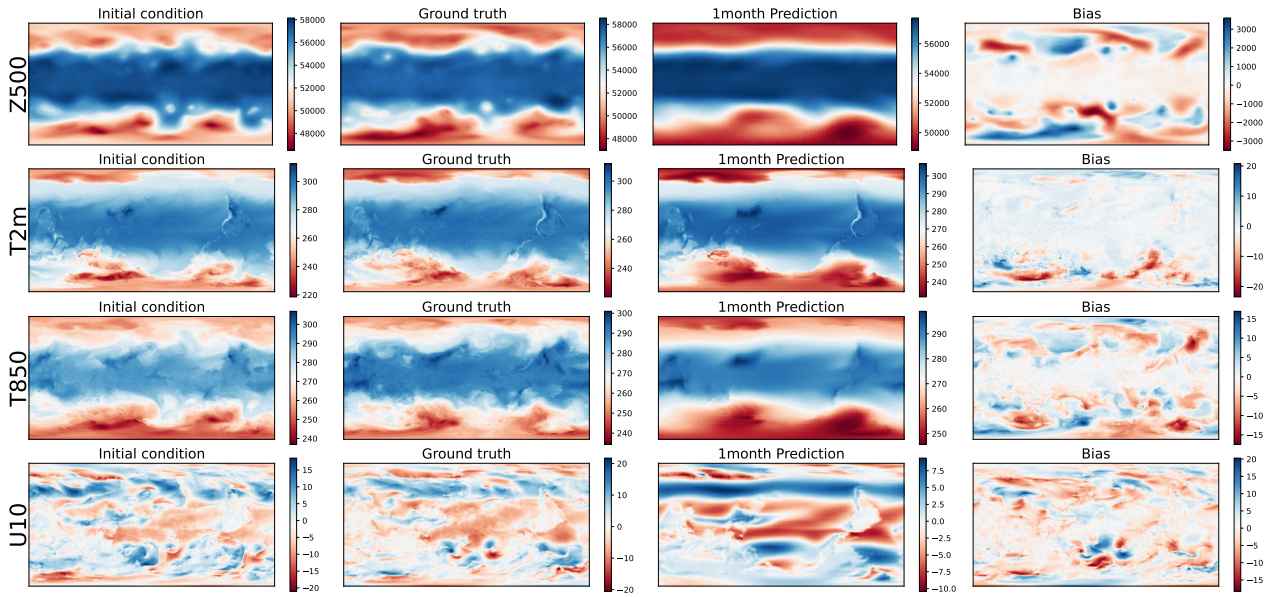**Figure 19:** Example forecasts from ClimaX at 2-week lead time compared to ground truth ERA5.



**Figure 20:** Example forecasts from ClimaX at 1-month lead time compared to ground truth ERA5.