
FENGWU: PUSHING THE SKILLFUL GLOBAL MEDIUM-RANGE WEATHER FORECAST BEYOND 10 DAYS LEAD

A PREPRINT

Kang Chen^{1,2,*} Tao Han^{1,*} Junchao Gong^{1,3,*} Lei Bai^{1,*,†} Fenghua Ling⁴ Jing-Jia Luo⁴

Xi Chen⁵ Leiming Ma⁶ Tianning Zhang¹ Rui Su¹ Yuanzheng Ci¹ Bin Li²

Xiaokang Yang³

Wanli Ouyang¹

* Equal Contributions, † Project Lead, bailei@pjlab.org.cn

¹ Shanghai Artificial Intelligence Laboratory

² University of Science and Technology of China

³ Shanghai Jiao Tong University

⁴ Nanjing University of Information Science and Technology

⁵ The Institute of Atmospheric Physics, Chinese Academy of Sciences

⁶ Shanghai Meteorological Bureau

April 7, 2023

ABSTRACT

We present FengWu, an advanced data-driven global medium-range weather forecast system based on Artificial Intelligence (AI). Different from existing data-driven weather forecast methods, FengWu solves the medium-range forecast problem from a multi-modal and multi-task perspective. Specifically, a deep learning architecture equipped with model-specific encoder-decoders and cross-modal fusion Transformer is elaborately designed, which is learned under the supervision of an uncertainty loss to balance the optimization of different predictors in a region-adaptive manner. Besides this, a replay buffer mechanism is introduced to improve medium-range forecast performance. With 39-year data training based on the ERA5 reanalysis, FengWu is able to accurately reproduce the atmospheric dynamics and predict the future land and atmosphere states at 37 vertical levels on a 0.25° latitude-longitude resolution. Hindcasts of 6-hourly weather in 2018 based on ERA5 demonstrate that FengWu performs better than GraphCast in predicting 80% of the 880 reported predictands, e.g., reducing the root mean square error (RMSE) of 10-day lead global z500 prediction from 733 to 651 m^2/s^2 . In addition, the inference cost of each iteration is merely 600ms on NVIDIA Tesla A100 hardware. The results suggest that FengWu can significantly improve the forecast skill and extend the skillful global medium-range weather forecast out to 10.75 days lead (with ACC of z500 > 0.6) for the first time.

Keywords Medium-range Weather Prediction · Deep Learning · Multi-modal Multi-task Learning · Transformer

1 Introduction

Understanding and predicting the Earth’s environment we live on, especially the atmosphere system, is a long-standing pursuit of human beings. Meteorological phenomena were recognized 3,000 years ago in the inscriptions on bones or tortoise shells of the Shang Dynasty in China (ca. sixteenth to eleventh century BCE) (Di, 2008). The atmosphere

system attracts more attention under global warming and the upsurge of extreme weather events (Lehmann et al., 2015; Rahmstorf and Coumou, 2011). Weather forecast, which involves analysis of past and present weather observations to predict future atmospheric conditions ranging from hours to days and even weeks, plays a critical role in routine decision-making for agriculture management, transportation, natural disasters (e.g., floods and storms) prevention, and green energy production, etc. For millennia, people have recognized the importance of weather forecast and have sought diverse ways to predict the weather. For example, the ancient Babylonians attempted to predict the weather based on sky observations (Taub, 2004) and the ancient Chinese developed instruments made of copper or wood to measure the wind speed.

Among various weather prediction tasks, global medium-range weather forecast, which targets predicting future global atmospheric conditions up to fourteen days ahead (Bengtsson, 1985), is arguably one of the most highly demanded tasks. It not only serves as the foundation for the deployed global ensemble forecast system (Gneiting and Raftery, 2005) to directly provide weather forecast services, but also provides background information and boundary conditions for regional numerical weather forecast systems (Mass and Kuo, 1998). Electronic computers were utilized in the 1950s for medium-range weather forecast by solving the partial differential equations (Bolin, 1955). Still, they were not widely available due to limitations in computing resources and data availability (Lynch, 2008). However, with the rapid development of Earth observation techniques (e.g., satellites) and High-Performance Computing (HPC) facilities, an increasing number of physical processes (e.g., radiation, thermodynamics, and fluid dynamics) could be simulated in a higher resolution with more accurate observations, leading to apparent forecast skills improvements from the 1980s.

Despite the significant breakthroughs achieved in the past decades, the performance of global medium-range weather forecast systems is still limited in forecast accuracy and extendibility due to large uncertainties in initial and boundary conditions, complicated non-linear physical processes, and heavy computation costs (Bauer et al., 2015). With the accumulation of massive weather observations and the maturity of deep learning techniques (e.g., large-scale training frameworks), researchers have commenced exploration on the possibility of AI-driven Numerical Weather Prediction (NWP) models. Specifically, Rasp et al. (Rasp et al., 2020) first introduced ResNet (He et al., 2016) to generate 5.625° latitude-longitude resolution of global weather prediction. Hu (Hu et al., 2023) utilized Recurrent Neural Network (RNN) architecture with variational loss to improve long-lead forecasts. While these attempts reveal the potential of data-driven methods in numerical weather prediction, they are limited in low-resolution data, leading to limited forecast applications. Recently, FourCastNet (Pathak et al., 2022), the first model producing 0.25° resolution forecasts, applies Vision Transformer (ViT) (Dosovitskiy et al., 2020) and Adaptive Fourier Neural Operators (AFNO) (Guibas et al., 2021) for efficient computation. Then, PanGu (Bi et al., 2022) acquires promising medium-range performance that surpasses ECMWF Integrated Forecasting System (IFS) in 0.25° resolution with a multi-timescale model combination strategy based on four 3D Earth-Specific Transformers. GraphCast (Lam et al., 2022) further boosts the AI methods' upper bound in NWP. It is more accurate in predicting 90% of the atmospheric variables compared with the ECMWF's deterministic operational forecasting system (IFS-HRES). In GraphCast, Graph neural network (GNN) is employed for medium-range global weather forecast, a 12-step autoregressive (AR) finetuning is adapted as the strategy for increasing the long-lead prediction accuracy, and mean squared error (MSE) loss used in GraphCast is elaborately weighted based on the pressure levels and weather variables.

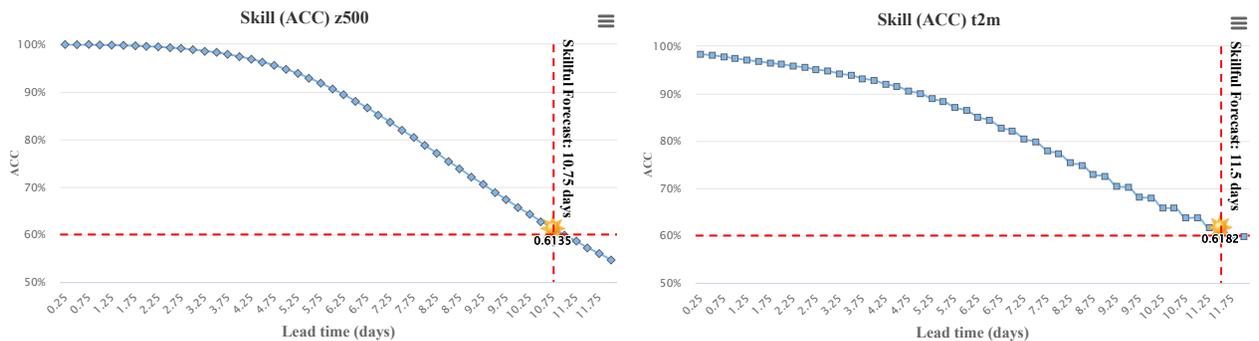


Figure 1: **The skillful forecast lead time of FengWu on z500 and t2m.** With the ACC > 0.6 as the criterion of a skillful weather prediction system, our findings are that FengWu can push the skillful forecast lead times to 10.75 days and 11.5 days for z500 and t2m, respectively.

In this paper, we propose an AI model to improve the medium-range weather forecast from a multi-modal and multi-task perspective. Specifically, our first proposition is to recognize the high-dimensional weather data, consisting of different atmospheric variables (e.g., temperature, winds, geopotential height, etc.) in different pressure levels, as distinct modalities, with each variable being treated as a modality; this is different from existing AI-based solutions that stack

all variables as single-modal input. This approach enables us to leverage existing multi-modal designs, e.g. processing the data of different atmosphere factors with modal-specific encoders, followed by a cross-modal Transformer to model the interactions among all atmosphere variables. The future states of each variable is derived separately after the cross-modal Transformer with modal-specific decoders. Our second proposition is to highlight that NWP is a multi-task regression problem based on the view that the prediction of each variable can be treated as a unique task. The prediction of some atmosphere variable will be more difficult than others, making the methods that treat different tasks equally unable to reach the global minimum point. Manually tuning the optimal weight for each task (variable) is prohibitively expensive and difficult. Treating NWP as a multi-task problem, we can leverage the uncertainty loss developed in the multi-task learning paradigm to the weather forecast domain, which considers the homoscedastic uncertainty of each variable. By learning to automatically scale weights of variable regression, introducing multi-task learning can improve learning effectiveness, leading to better weather prediction accuracy.

Another issue in AI-driven global medium-range weather forecast is generating long-lead predictions (e.g., 10 days' predictions), which is hard to be directly optimized due to the extremely huge volume of global weather data (e.g., 2.3 GB for each time step in the ERA5 pressure level dataset with float32 format) even with the most advanced GPU devices. Previous works tackled this problem either by autoregressively fine-tuning with two (Pathak et al., 2022) or multiple steps (Lam et al., 2022) powered by elaborate engineering techniques, or training separate models for different steps (Bi et al., 2022). Despite they are demonstrated to be effective, these solutions can be computationally expensive and memory-intensive, making them challenging to implement when computation resources are limited. To solve the long-lead prediction issue, we propose the Replay Buffer mechanism, which is inspired by the reinforcement learning study (Schaul et al., 2015). The replay buffer stores the predicted results from previous optimization iterations and uses them as the current model's input, which mimics the intermediate input error during the auto-regressive inference stage. This design boosts the long-lead forecast quality with efficient computation and memory.

Based on the above perspectives and designs, we develop FengWu¹, an advanced weather forecast system for global medium-range weather predictions. By training with the high-resolution (i.e., 0.25° latitude-longitude resolution) ERA5 dataset over the past 39 years, FengWu is able to accurately emulate the atmospheric dynamics and predict the future land and atmosphere status of 37 levels. Hindcasts of 6-hourly weather in the year 2018 indicate that FengWu achieves the best forecast skills among all the released data-driven forecast systems. Specifically, FengWu has higher accuracy than GraphCast on 80% of the 880 reported predictands. The forecast skill improvements make the skillful forecast lead time of FengWu reach 10.75 days (ACC of z500 > 0.6 Bauer et al. (2015)) for the first time by an AI-based approach.

2 Preliminary

2.1 Dataset

ERA5 (Hersbach et al., 2020) is a global atmospheric reanalysis dataset produced by the European Centre for Medium-Range Weather Forecasts (ECMWF). It provides comprehensive information about the Earth's climate and weather conditions, covering the period from 1940 to the present. ERA5 provides a wide range of variables such as temperature, humidity, precipitation, wind speed and direction, mean sea level pressure, and many others. The data is available at a high spatial resolution of 0.25° latitude-longitude resolution and 37 vertical pressure levels, ranging from 1000 hPa to 1 hPa, which makes it suitable for a wide range of applications, including climate research, weather forecast, and environmental monitoring.

In this study, FengWu simulates 5 atmospheric variables (each with 37 pressure levels) and four surface variables, a total of 189 predictands. Specifically, the atmospheric variables are geopotential (z), relative humidity (r), zonal component of wind (u), meridional component of wind (v), and air temperature (t), whose 37 sub-variables at different vertical level are presented by abbreviating their short name and pressure levels (e.g., z500 denotes the geopotential height at a pressure level of 500 hPa). And the four surface variables are 2-meter temperature (t2m), 10-meter u wind component (u10), 10-meter v wind component (v10), and mean sea level pressure (msl).

For consistency, we follow the validation strategies demonstrated by GraphCast, i.e., the data from 1979-2015 is used for training, 2016-2017 for validation, and 2018 for testing. In addition, we also leverage the 6-hourly sampled data (T00, T06, T12, T18) instead of the hourly ERA5 dataset for training.

¹The name of "FengWu" comes from the ancient Chinese anemometer used from the Han dynasty, considered as the earliest prototype for measuring wind speed and orientation.

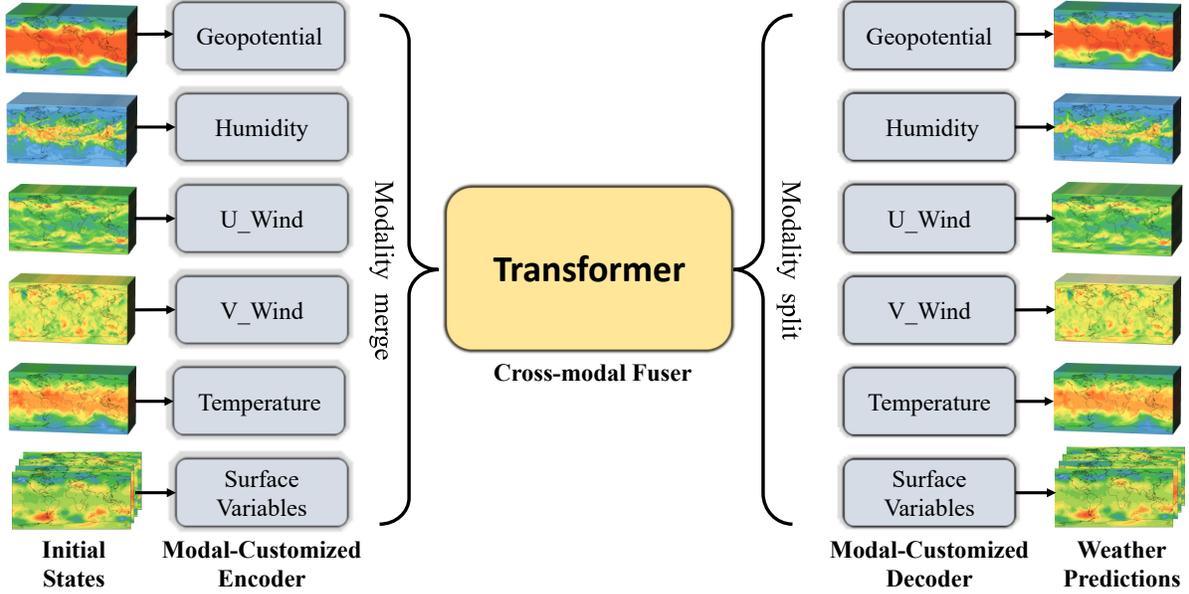


Figure 2: Overview of FengWu’s architecture. FengWu first treats the multiple weather factors as different modalities and extracts their feature embeddings independently. And then a transformer-based network is utilized to fuse and pass messages among different modalities. Finally, the high-level feature representation is used to get the predictors via the modal-customized decoder.

2.2 Problem Formulation

The objective of FengWu is to utilize AI techniques to build the most skillful high-resolution medium-range global weather forecast system, which predicts the future 14-day global atmosphere states based on the current atmosphere conditions. Formally, we denote the weather state at time slot i as a high dimension tensor $X^i \in \mathbb{R}^{C \times W \times H}$, where C denotes the number of atmosphere variables considered in this work, W and H are the width and height. When mapping the continuous atmosphere fluid to a 2D spatial plane with 0.25° latitude-longitude resolution, we have $C = 189$, $W = 721$, and $H = 1440$. FengWu aims at generating 14-day lead time forecasts $\{\hat{X}^{i+1}, \hat{X}^{i+2}, \dots, \hat{X}^{i+56}\}$ with a time-interval of six hours,

$$\{\hat{X}^{i+1}, \hat{X}^{i+2}, \dots, \hat{X}^{i+56}\} = \text{FengWu}(X^i), \quad (1)$$

where $\hat{X}^{i+\tau}$ is the prediction of the weather state at time slot $i + \tau$. However, it is hard to directly learn a function for Eq. 1 due to the extremely large size of the global atmosphere data. Following the practice of atmosphere simulation, FengWu targets on learning a function to predict the data of the next step, which could then generate multi-step predictions in an auto-regressive manner, i.e.,

$$\hat{X}^{i+1} = \text{FengWu}(X^i), \hat{X}^{i+2} = \text{FengWu}(\hat{X}^{i+1}), \dots, \hat{X}^{i+56} = \text{FengWu}(\hat{X}^{i+55}) \quad (2)$$

3 Method

In this section, we introduce the details of FengWu, which includes three main components: 1) the network Architecture including Transformer-based modal-customized encoder-decoders and cross-modal fuser, 2) the uncertainty loss for multi-task optimization, and 3) the replay buffer for long-lead predictions.

3.1 Network Architecture

FengWu considers weather variables as different modalities of the atmosphere state and employs an "encode-fuse-decode" structure as illustrated in Figure 2. The *Modal-Customized Encoder* encodes multi-modal features, which are fused by a Transformer-based *Cross-modal Fuser* to get the joint representations. The *Modal-Customized Decoder* then separately predicts the weather variables from the joint representations.

Modal-Customized Encoder. We design Modal-Customized Encoders to extract features independently. Each weather state X^i with shape (C, W, H) is sliced into earth surface state X_s^i , geopotential state X_z^i , humidity state X_q^i , the eastward component of the wind state X_u^i , the northward component of wind state X_v^i , and temperature state X_t^i whose shapes are, respectively, (C_s, W, H) , (C_z, W, H) , (C_q, W, H) , (C_u, W, H) , (C_v, W, H) , and (C_t, W, H) . To obtain features \tilde{X}_m for $m \in \{s, z, q, u, v, t\}$ separately, a transformer-based encoder $f_{en,m}(X_m|\theta_{en,m})$ with encoder parameters $\theta_{en,m}$ is used for its corresponding state X_m . The output of the encoder is denoted as Z_m for the modality m , where

$$Z_m = f_{en,m}(X_m|\theta_{en,m}). \quad (3)$$

Cross-modal Fuser. The output of the encoder Z_m for $m \in \{s, z, q, u, v, t\}$ are concatenated to obtain the fused features as follows:

$$Z = \text{concat}(Z_s, Z_z, Z_q, Z_u, Z_v, Z_t), \quad (4)$$

where *concat* denotes feature concatenation along the feature channel dimension. The fused features are then fed into a transformer for fusing their information and extracting the fused features \tilde{Z} .

Modal-Customized Decoder. In Multi-task Decoder, tokens generated by Multi-modal Feature Fuser are used for predicting the mean and variance of atmosphere variables. Separate modality decoders $f_{de,m}(\tilde{Z}|\theta_{de,m})$ for $m \in \{s, z, q, u, v, t\}$ are designed to predict the future state of corresponding modalities, where $\theta_{de,m}$ denotes the parameters of the decoder $f_{de,m}$ for modality m . The decoder network is similar to the encoder.

3.2 Uncertainty Loss for Multi-task Optimization

In this paper, weather forecast learning is regarded as multi-task learning. Previous studies in the multi-task learning paradigm show that assigning different weights between different tasks is helpful for learning representations. This observation is consistent with the practice in GraphCast (Lam et al., 2022), which assigns manually designed weights for different variables and pressure levels as hyperparameters. While it is demonstrated to be effective, manually determining approximate weights for variables would be arduous and sub-optimal.

To solve the multi-task optimization issue more accurately and elegantly, we introduce the uncertainty loss to automatically learn weights for weather forecasting. Specifically, FengWu is defined as a probabilistic model that predicts the parameters $\hat{\mu}^{i+1}, \hat{\sigma}^{i+1}$ of a Gaussian distribution:

$$\hat{\mu}^{i+1}, \hat{\sigma}^{i+1} = \text{FengWu}(X^i) \quad (5)$$

where $\hat{\mu}^{i+1}$ and $\hat{\sigma}^{i+1}$ are predicted mean and variance of predictands X_{i+1} . And the probability of atmosphere variables can be calculated by the mean and variance:

$$p(x_{c,w,h}^{i+1} | \hat{\mu}^{i+1}, \hat{\sigma}^{i+1}) = \mathcal{N}\left(\hat{\mu}_{c,w,h}^{i+1}, \hat{\sigma}_{c,w,h}^{i+1}\right) \quad (6)$$

In Eq. 6, each element $x_{c,w,h}^{i+1}$ in X^{i+1} with subscript (c, w, h) follows an independent univariate Gaussian distribution $\mathcal{N}\left(\hat{\mu}_{c,w,h}^{i+1}, \hat{\sigma}_{c,w,h}^{i+1}\right)$, where $c = 1, \dots, 189$ denotes the index for the channel, i.e. different pressure levels and weather variables, e.g. temperature. w and h respectively denote the latitude grid and longitude grid. We adopt maximum likelihood estimation to allocate weights for different tasks(variables). Because we employ the likelihood as the minimization objective, the loss of each variable c at location (w, h) is automatically weighted by homoscedastic uncertainty. The uncertainty loss provides an approach to tradeoff the weights between variables, pressure levels, and locations without an expensive manual search, which is efficient, particularly in large weather forecast models.

3.3 Replay Buffer for Long-lead Predictions

Trained as a single-step predictor, directly using FengWu for medium-range prediction will lead to inferior long-lead forecast performance due to the accumulation of errors in the AutoRegressive (AR) inference process. GraphCast effectively eases this problem by adding an autoregressive training stage, which gradually increases the number of autoregressive steps in the training from 2 to 12. However, this approach encounters two challenges. Firstly, the intermediate predictions generated during multiple forwards can be wasteful, as they are discarded after being optimized and are not reused for other autoregressive steps, which reduces their efficiency and slows down the training process. Secondly, the memory required to store and process the gradient in the constantly increasing autoregressive steps can become excessive, limiting the maximum AR steps that can be processed.

To tackle the problems mentioned above, this work proposes a replay buffer mechanism. Specifically, a set $\mathcal{B} = \{\hat{X}_j^{i+\tau}\}_{j=0}^N$ containing N predictions are denoted to represent the data in the buffer. Initially, the replay buffer first pushes a certain number of first-step predictions in the initial stage. In the next stage, FengWu learns from both the original dataset and the replay buffer, e.g., sampling data from both the original dataset and the replay buffer. Accordingly, the predicted results, either taking input from the original data or the replay buffer, are treated as the intermediate predictions and saved to the data buffer, resulting in diverse finetuning frequencies between different autoregressive steps. The last element in the replay buffer will be popped out if the queue is full for each data collection. The replay buffer plays a critical role in enabling our system to perform long-lead autoregressive forecasts by collecting and reusing intermediate predictions, thus enforcing the FengWu to take the accumulated autoregressive estimation errors into consideration during training. This online learning strategy is particularly valuable in situations where the devices or framework does not support long-lead AR training. Specifically, a training sample at time step $i + 1$ stored in the replay buffer is used as the input of FengWu for prediction at time step $i + 2$ and its predicted results at time step $i + 2$ are stored in the replay buffer, and the predicted results at time step $i + 2$ in the replay buffer will be used as the input in the latter training iteration. Therefore, the replay buffer helps the training stage to simulate the long-lead autoregressive forecast during the inference in Eq. 2. In addition to its role in enabling long-lead AR forecast, the replay buffer also offers the benefit of reducing GPU memory usage by storing data on the CPU. Overall, the replay buffer enhances the efficiency and effectiveness of the long-lead AR learning process.

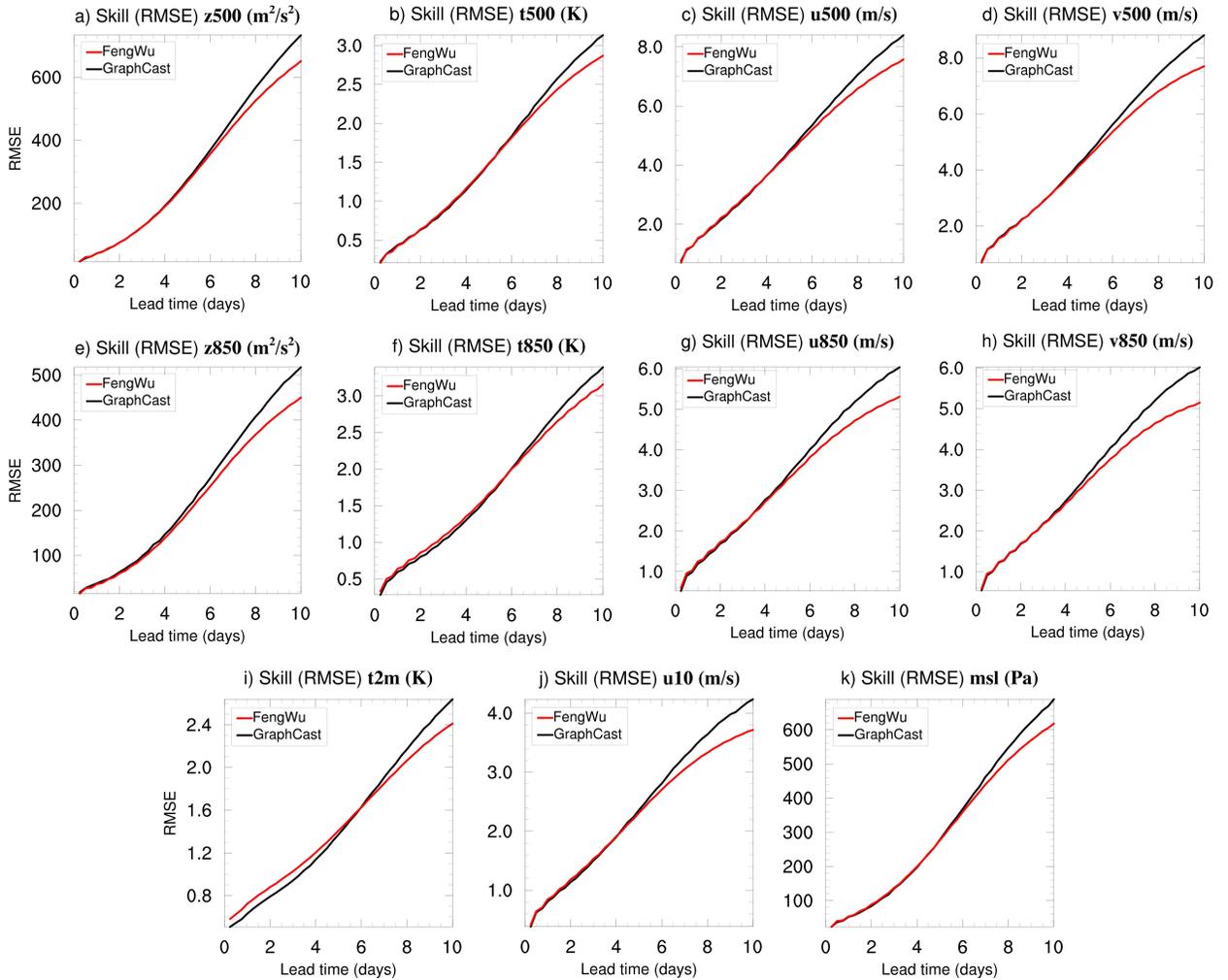


Figure 3: **Latitude-weighted RMSE skill of FengWu (red lines) and GraphCast (black lines) predicting the weather in 2018 (Lower RMSE is better).** The x-axis in each sub-figure represents lead time, at a 6-hour interval over a 10-day lead time. The y-axis represents the latitude-weighted RMSE defined in Eq. 7.

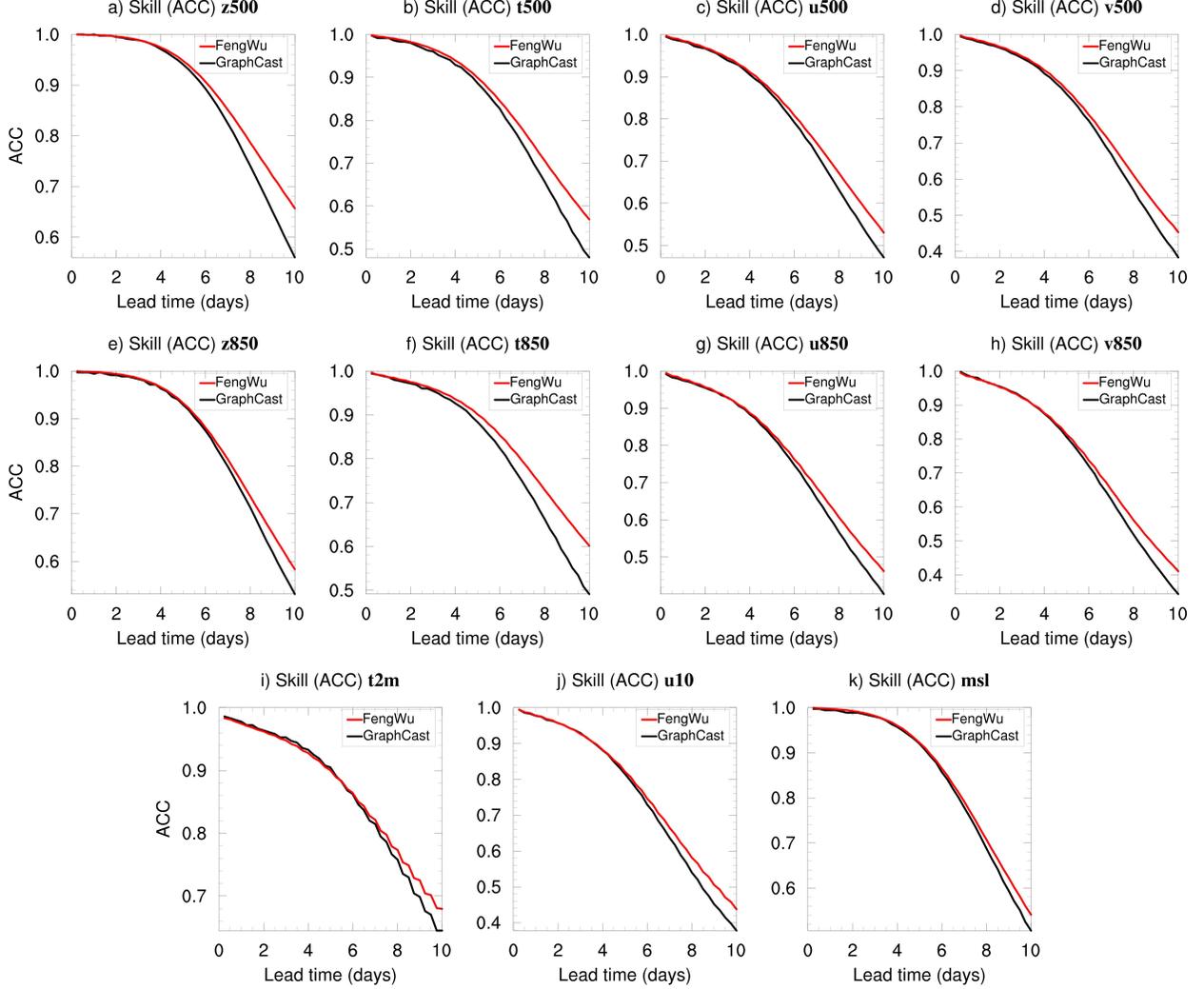


Figure 4: ACC skill of FengWu and GraphCast predicting the weather in 2018 (Higher ACC is better). The x-axis in each sub-figure represents lead time, at a 6-hour interval over a 10-day lead time. The y-axis represents the ACC defined in Eq. 8.

4 Results

4.1 Evaluation Strategies

For consistency, we follow the evaluation protocols implemented in the work by GraphCast, which includes the same evaluation metrics, dataset splitting, and lead time of forecast. With 00z and 12z as the initial weather states for each day, we compare the performance of FengWu and GraphCast for a 10-day forecast using the commonly used RMSE and ACC metrics based on the test set.

RMSE represents the latitude-Weighted Root Mean Square Error, which is a statistical metric widely used in geospatial analysis and climate science to assess the accuracy of a model’s predictions or estimates of temperature, precipitation, or other meteorological variables across different latitudes. Given the prediction result $\hat{x}_{c,w,h}^{i+\tau}$ and its target (ground truth) $x_{c,w,h}^{i+\tau}$, the RMSE is defined as follows:

$$\text{RMSE}(c, \tau) = \frac{1}{T} \sum_{i=1}^T \sqrt{\frac{1}{W \cdot H} \sum_{w=1}^W \sum_{h=1}^H W \cdot \frac{\cos(\alpha_{w,h})}{\sum_{w'=1}^W \cos(\alpha_{w',h})} (x_{c,w,h}^{i+\tau} - \hat{x}_{c,w,h}^{i+\tau})^2}, \quad (7)$$

where c denotes the index for channels that are either the surface variable or the atmosphere variable at a certain pressure level. w and h respectively denote the indices for each grid along the latitude and longitude indices. $\alpha_{w,h}$ is the latitude of point (w, h) . T is the total number of valid test time slots.

ACC is the Latitude-weighted Anomaly Correlation Coefficient that evaluates the performance of dynamical models by comparing their predictions of anomalies (departures from the long-term averaged climatology) to observed anomalies. ACC is similar to the standard Anomaly Correlation Coefficient, but also with a latitude weighting factor applied to account for the varying area represented by different latitudes on a spherical Earth,

$$\text{ACC}(c, \tau) = \frac{1}{T} \sum_{i=1}^T \frac{\sum_{w,h} W \cdot \frac{\cos(\alpha_{w,h})}{\sum_{w'=1}^W \cos(\alpha_{w',h})} (x_{c,w,h}^{i+\tau} - C_{c,w,h}^{i+\tau})(\hat{x}_{c,w,h}^{i+\tau} - C_{c,w,h}^{i+\tau})}{\sqrt{\sum_{w,h} W \cdot \frac{\cos(\alpha_{w,h})}{\sum_{w'=1}^W \cos(\alpha_{w',h})} (x_{c,w,h}^{i+\tau} - C_{c,w,h}^{i+\tau})^2 \sum_{w,h} W \cdot \frac{\cos(\alpha_{w,h})}{\sum_{w'=1}^W \cos(\alpha_{w',h})} (\hat{x}_{c,w,h}^{i+\tau} - C_{c,w,h}^{i+\tau})^2}}, \quad (8)$$

where $C_{c,w,h}^{i+\tau}$ is the climatological mean over the day-of-year containing the validity time $i + \tau$ for a given weather variable c at longitude w and latitude h . It is averaged from the years 1993 to 2016 with the ERA5 data on a daily basis, which is consistent with the approach taken by GraphCast. Note that if $C_{c,w,h}^{i+\tau}$ is calculated hourly, the resulting ACC metric may be significantly higher than the values obtained using the daily climate mean. This is because the ACC metric is sensitive to outliers and can be influenced by various factors, such as the choice of the time period and spatial resolution of the data.

4.2 Quantitative Skill Evaluation

Figures 3 and 4 illustrate the predictive performance comparison of FengWu (red lines) and GraphCast (black lines) in terms of RMSE and ACC, respectively. The analysis is conducted over 880 targets at 6-hour intervals, with a lead time of 10 days. The results show that FengWu has both lower RMSE and higher ACC than GraphCast on 80% of the targets analyzed. FengWu demonstrates a comparable level of forecasting accuracy as GraphCast for predicted variables within a lead time range of 1~5 days, except for t2m. In particular, as the lead time increases, significant improvement with FengWu is observed for all predicted variables, demonstrating FengWu’s remarkable ability for long-lead weather forecasting.

4.3 Qualitative Prediction Evaluation

We visualize the predicted results of FengWu at lead days 3, 5, and 10 for two variables, i.e., z500 (geopotential at the pressure level of 500 hPa) and t850 (the temperature at the pressure level of 850 hPa), and compare the predictions with the ERA5. In Figure 5 and Figure 6, the top two rows show the sequences of states from FengWu and ERA5, and the third row shows the absolute value of the error from FengWu to ERA5. In both visualizations, FengWu has outcomes close to ERA5 on the third day. As the forecast step increases, the absolute error increases and diffuses to adjacent areas. These visualizations validate FengWu’s ability to estimate weather states approximating the real data.

4.4 Effects of the Replay Buffer

As mentioned in Section 3.3, the primary goal of the replay buffer is to reduce the accumulated error in the long-term prediction caused by the autoregressive estimation paradigm. To evaluate the effectiveness of the replay buffer in boosting forecasting skills, we compare the performance of FengWu with and without the replay buffer mechanism. As shown in Figure 7, the system’s forecast performance decreases significantly as lead time increasing when the proposed replay buffer is removed, indicating that it is a crucial component in improving the accuracy of long-lead weather predictions.

4.5 Computation Cost

Training Cost: FengWu is developed with Pytorch and trained with 32 Nvidia A100 GPUs in a cluster for 17 days in total. Compared with GraphCast trained with 32 Cloud TPU v4 for 21 days, we consider the training time of FengWu is only 47% to 67% of GraphCast (according to Jouppe et al. (2023), TPU v4 is 1.2 to 1.7 times faster than Nvidia A100).

Inference Cost: We evaluate the inference speed of FengWu on an NVIDIA Tesla-A100 GPU, which indicates that FengWu costs less than 30 seconds to generate all forecasts in the following 10 days with a six-hour interval. With a peak power consumption of 0.4KW for an A100 (Choquette et al., 2021), a 10-day inference by FengWu consumes

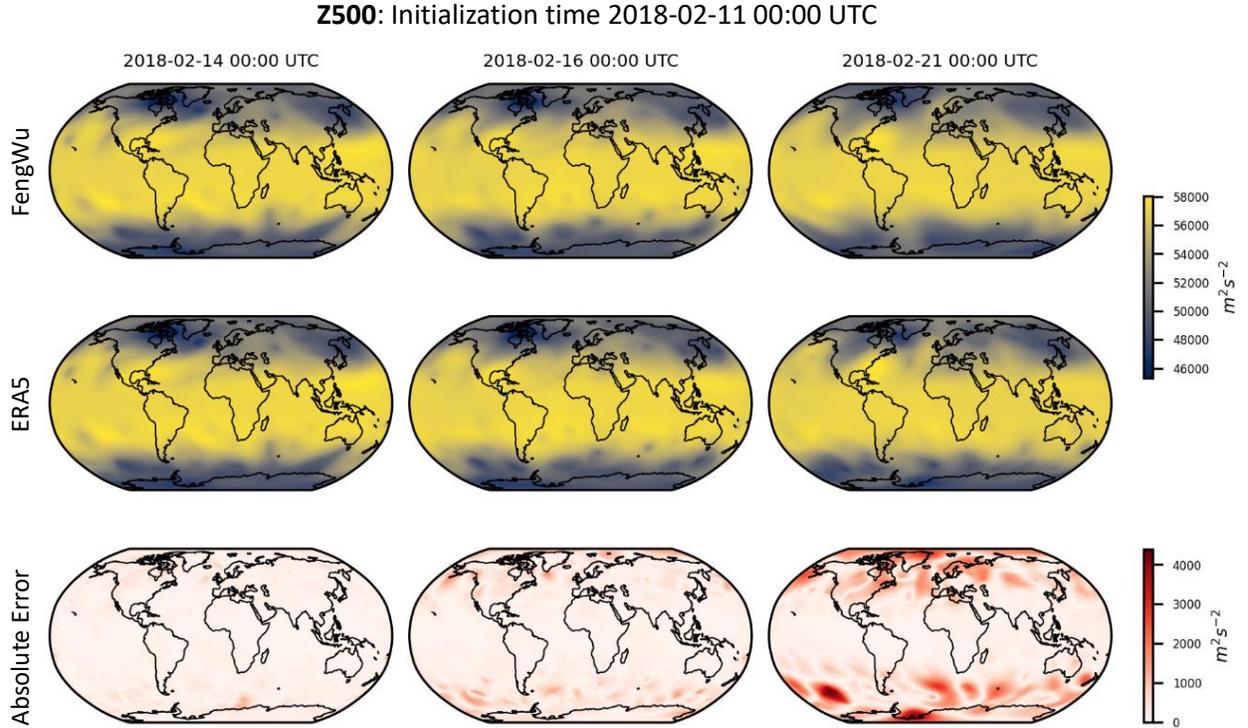


Figure 5: **Forecast images and absolute error for z500.** Figures of z500 on days 3, 5, and 10 are presented with initialization time at 2018-02-11 00:00 UTC. The subtitles at the top of the columns indicate the dates of prediction. The first row and second row show FengWu and ERA5 ground truth, respectively. Row 3 shows the absolute error between FengWu and ERA5.

roughly 12kJ energy, while the consumption of a single member of the IFS model is estimated to be about 26.6MJ², approximately 2000 times higher than FengWu.

5 Conclusions and Discussions

In this paper, we introduce FengWu, an advanced AI-based weather forecasting system, which has three technical contributions. First, we propose to solve the problem of global medium-range weather forecasting as a multi-modal multi-task learning problem and introduce corresponding techniques for it, including multi-modal network architecture and uncertainty-based multi-task loss. Second, we propose a replay buffer mechanism, which could improve the long-term forecasting performance under the autoregressive inference setting with limited devices. Third, FengWu achieves top performance among all existing AI-based methods and extends the skillful global medium-range weather forecast lead time to 10.75 days (ACC of z500 > 60%). It also has higher accuracy than GraphCast, the current state-of-art AI-based weather forecasting model, on 80% reported prediction targets.

As the initial fields of uneven quality are applied in physics-based and AI-based models, the fairness of model comparisons is supposed to be discussed. Specifically, the precision of weather forecasting systems greatly depends on the initial fields because the weather system is chaotic, and the chaotic system obtains precise estimations of its future with accurate initial values. AI methods including FengWu and GraphCast (Lam et al., 2022) forecast weather with initial states from ERA5, while physics systems like IFS-HRES use their own initial analysis. To achieve timely forecasts, IFS-HRES produces such analysis fields with observations accessible at the start moment of prediction. For example, if observations from some satellites take several hours to be synchronized, they might be excluded from data assimilation for initial analysis states. Compared with IFS-HRES, ERA5 provides a more complete and accurate picture of initial fields, as the 5 days delay in ERA5’s analysis data allows all observations to be blended. As a consequence, AI

²82 minutes is required by the IFS model in a 15-day, 51-member ensemble forecast applying 18km resolution grid on 1530 Cray XC40 nodes with dual-socket Intel Haswell processors (Bauer et al., 2020). A dual-socket Intel Haswell node draws a Thermal Design Power (TDP) of 270 Watts (Pathak et al., 2022). The 10-day forecasts for a single IFS member roughly take 98,400 node seconds.

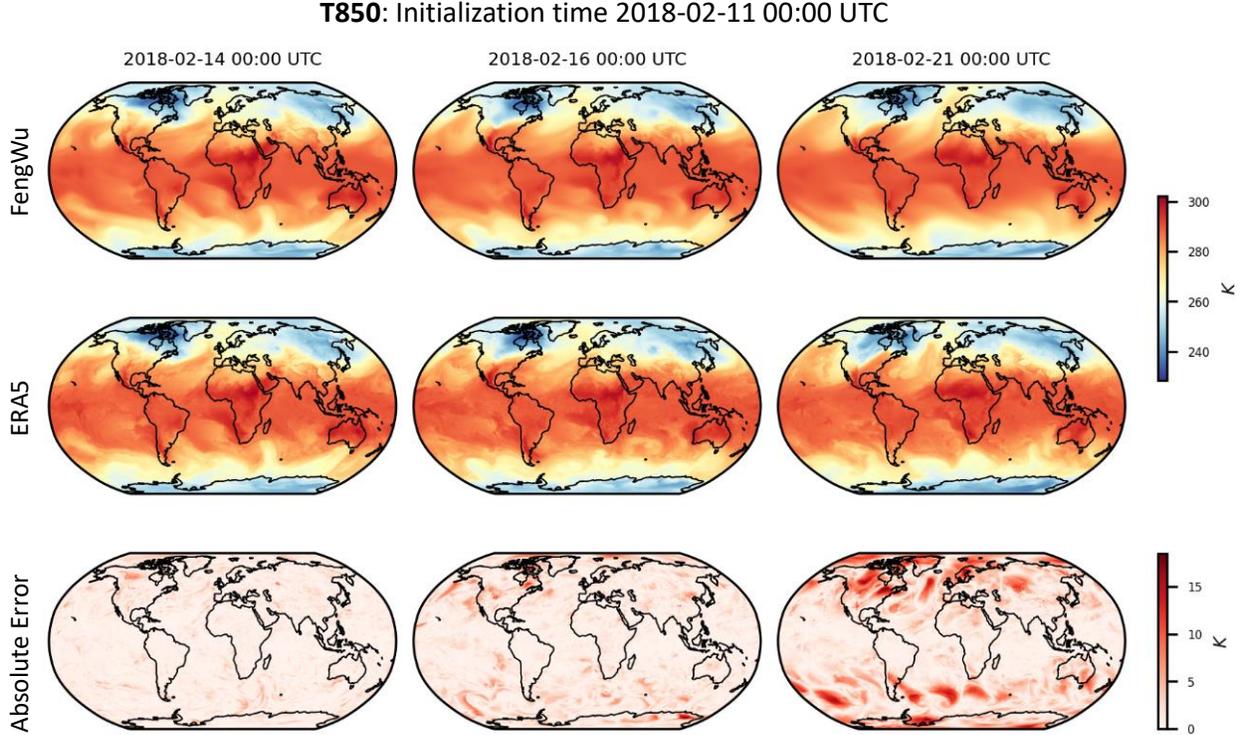


Figure 6: **Forecast images and absolute error for T850.** The plot demonstrates ERA5 ground truth and FengWu’s prediction for T850, with forecast initialization 2018-02-11 00:00. Other settings are similar to Figure 5.

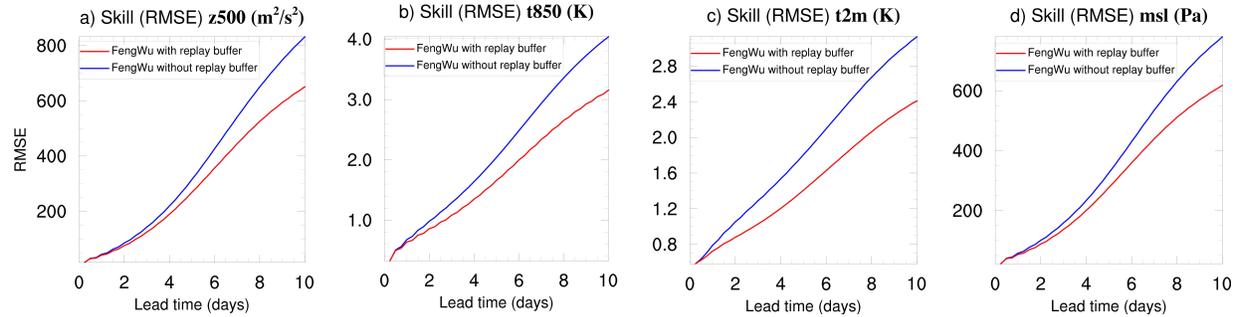


Figure 7: **Effects of the autoregressive training with the proposed replay buffer mechanism (red lines) or without it (blue lines).** The x-axis in each sub-figure represents the lead time at 6-hour steps over 10 days. The y-axis represents RMSE defined in Eq. 7 (lower is better).

methods occupy a favorable position, for ERA5’s premium initial states, in comparisons with IFS-HRES, while the comparison between AI methods, i.e., FengWu and GraphCast, is reasonable because their initial fields are both ERA5 reanalysis data.

Acknowledgements

We acknowledge the use of the ERA5 dataset on both pressure levels and single level provided by the European Centre for Medium-Range Weather Forecasts (ECMWF). Without their great efforts in collecting, archiving, and disseminating the data, this study would not be possible.

We acknowledge the Research Support, IT, and Infrastructure team based in the Shanghai AI Laboratory for their provision of computation resources and network support. F Ling and J-J Luo are supported by National Key Research and Development Program of China (No. 2020YFA0608000). This acknowledgment extends particularly to the

individuals on the team, namely Prof. Yu Qiao, Liang Liu, Qihong Liao, Jiamin Ge, Jing Zou, Jingwen Li, and Xingpu Li.

We would also like to express our appreciation to Prof. Yang Wang from the University of Science and Technology of China, Prof. Tao Chen from Fudan University, Prof. Hongsheng Li from The Chinese University of Hong Kong, Dr. Jiajun Deng from the University of Sydney, Dr. Tong He from Shanghai AI Laboratory, and Mr. Peng Ye for their help and valuable discussions during the conduction of this research, which have significantly enhanced the quality of this work. We are grateful for their contributions and acknowledge their role in the development of this work.

References

- Li Di. *Meteorology in China*, pages 1662–1664. Springer Netherlands, Dordrecht, 2008. ISBN 978-1-4020-4425-0. doi:[10.1007/978-1-4020-4425-0_8787](https://doi.org/10.1007/978-1-4020-4425-0_8787). URL https://doi.org/10.1007/978-1-4020-4425-0_8787.
- Jascha Lehmann, Dim Coumou, and Katja Frieler. Increased record-breaking precipitation events under global warming. *Climatic Change*, 132:501–515, 2015.
- Stefan Rahmstorf and Dim Coumou. Increase of extreme events in a warming world. *Proceedings of the National Academy of Sciences*, 108(44):17905–17909, 2011.
- Liba Taub. *Ancient meteorology*. Routledge, 2004.
- Lennart Bengtsson. Medium-range forecasting at the ecmwf. In *Advances in Geophysics*, volume 28, pages 3–54. Elsevier, 1985.
- Tilmann Gneiting and Adrian E Raftery. Weather forecasting with ensemble methods. *Science*, 310(5746):248–249, 2005.
- Clifford F Mass and Ying-Hwa Kuo. Regional real-time numerical weather prediction: Current status and future potential. *Bulletin of the American Meteorological Society*, 79(2):253–264, 1998.
- Bert Bolin. Numerical forecasting with the barotropic model 1. *Tellus*, 7(1):27–49, 1955.
- Peter Lynch. The origins of computer weather prediction and climate modeling. *Journal of computational physics*, 227(7):3431–3444, 2008.
- Peter Bauer, Alan Thorpe, and Gilbert Brunet. The quiet revolution of numerical weather prediction. *Nature*, 525(7567):47–55, 2015.
- Stephan Rasp, Peter D Dueben, Sebastian Scher, Jonathan A Weyn, Soukayna Mouatadid, and Nils Thuerey. Weather-bench: a benchmark data set for data-driven weather forecasting. *Journal of Advances in Modeling Earth Systems*, 12(11):e2020MS002203, 2020.
- 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*, pages 770–778, 2016.
- Yuan Hu, Lei Chen, Zhibin Wang, and Hao Li. Swinvrnn: A data-driven ensemble forecasting model via learned distribution perturbation. *Journal of Advances in Modeling Earth Systems*, 15(2):e2022MS003211, 2023.
- 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. *arXiv preprint arXiv:2202.11214*, 2022.
- 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. *arXiv preprint arXiv:2010.11929*, 2020.
- John Guibas, Morteza Mardani, Zongyi Li, Andrew Tao, Anima Anandkumar, and Bryan Catanzaro. Adaptive fourier neural operators: Efficient token mixers for transformers. *arXiv preprint arXiv:2111.13587*, 2021.
- 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. *arXiv preprint arXiv:2211.02556*, 2022.
- Remi Lam, Alvaro Sanchez-Gonzalez, Matthew Willson, Peter Wirsberger, Meire Fortunato, Alexander Pritzel, Suman Ravuri, Timo Ewalds, Ferran Alet, Zach Eaton-Rosen, et al. Graphcast: Learning skillful medium-range global weather forecasting. *arXiv preprint arXiv:2212.12794*, 2022.
- Tom Schaul, John Quan, Ioannis Antonoglou, and David Silver. Prioritized experience replay. *arXiv preprint arXiv:1511.05952*, 2015.

- 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. *Quarterly Journal of the Royal Meteorological Society*, 146(730):1999–2049, 2020.
- Norman P. Jouppi, George Kurian, Sheng Li, Peter Ma, Rahul Nagarajan, Lifeng Nai, Nishant Patil, Suvinay Subramanian, Andy Swing, Brian Towles, Cliff Young, Xiang Zhou, Zongwei Zhou, and David Patterson. Tpu v4: An optically reconfigurable supercomputer for machine learning with hardware support for embeddings. *arXiv preprint arXiv:2304.01433*, 2023.
- Jack Choquette, Wishwesh Gandhi, Olivier Giroux, Nick Stam, and Ronny Krashinsky. Nvidia a100 tensor core gpu: Performance and innovation. *IEEE Micro*, 41(2):29–35, 2021.
- 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.