

1
2

Advancing Parsimonious Deep Learning Weather Prediction using the HEALPix Mesh

3
4
5

Matthias Karlbauer¹, Nathaniel Cresswell-Clay², Dale R. Durran²,
Raul A. Moreno², Thorsten Kurth³, Boris Bonev³, Noah Brenowitz⁴, and
Martin V. Butz¹

6
7
8
9
10

¹Neuro-Cognitive Modeling Group, Department of Computer Science, University of Tübingen, Tübingen,
Germany
²Department of Atmospheric Sciences, University of Washington, Seattle, WA, USA
³NVIDIA Switzerland AG, Zürich, Switzerland
⁴NVIDIA Corporation, Seattle, USA

11
12
13
14
15
16
17

Key Points:

- A U-Net is refined to forecast seven atmospheric variables on global scale, falling behind the state-of-the-art by only one day.
- Forecasts are generated on the HEALPix mesh, facilitating the development of location invariant convolution kernels.
- Without converging to climatology, the model produces stable and realistic states of the atmosphere in 365-days rollouts.

Abstract

We present a parsimonious deep learning weather prediction model on the Hierarchical Equal Area isoLatitude Pixelization (HEALPix) to forecast seven atmospheric variables for arbitrarily long lead times on a global approximately 110 km mesh at 3h time resolution. In comparison to state-of-the-art machine learning weather forecast models, such as Pangu-Weather and GraphCast, our DLWP-HPX model uses coarser resolution and far fewer prognostic variables. Yet, at one-week lead times its skill is only about one day behind the state-of-the-art numerical weather prediction model from the European Centre for Medium-Range Weather Forecasts. We report successive forecast improvements resulting from model design and data-related decisions, such as switching from the cubed sphere to the HEALPix mesh, inverting the channel depth of the U-Net, and introducing gated recurrent units (GRU) on each level of the U-Net hierarchy. The consistent east-west orientation of all cells on the HEALPix mesh facilitates the development of location-invariant convolution kernels that are successfully applied to propagate global weather patterns across our planet. Without any loss of spectral power after two days, the model can be unrolled autoregressively for hundreds of steps into the future to generate stable and realistic states of the atmosphere that respect seasonal trends, as showcased in one-year simulations. Our parsimonious DLWP-HPX model is research-friendly and potentially well-suited for sub-seasonal and seasonal forecasting.

Plain Language Summary

Weather forecasting is traditionally realized by numerical weather prediction models that solve physical equations to simulate the progression of the atmosphere. Numerical methods are compute intense and their performance is increasingly challenged by less compute demanding but still highly sophisticated machine learning approaches. Yet, a downside of these new models is their reliability: They are not guaranteed to generate physically plausible states, which often prevents them from generating stable and realistic forecasts beyond two weeks into the future. Here, a parsimonious machine learning model is developed to forecast just seven variables of the atmosphere (compared to more than 800 in numerical models and 67 or 218 in competitive machine learning models) over an entire year. Despite the small number of variables, our model generates forecasts that only fall behind expensive state-of-the-art predictions by a single day. That is, our error in a seven-days forecast matches that of a state-of-the-art forecast at day eight. Advancing weather forecasts with research friendly and parsimonious machine learning models beyond two weeks promises to extend horizons for planning in various fields that impact environment, economy, and society.

1 Introduction

Four years ago, Weyn et al. (2019) posed the question “Can machines learn to predict the weather?” and demonstrated that data driven convolutional neural networks can forecast the evolution of the 500 hPa surface much better than the alternative dynamical model, the barotropic vorticity equation, which was used in the first numerical weather prediction (NWP) model (Charney et al., 1950). An extremely rapid evolution of deep learning weather prediction (DLWP) models followed, culminating in the recent Pangu-Weather (Bi et al., 2023) and GraphCast models (Lam et al., 2022), which outperform the deterministic forecast from the state-of-the-art Integrated Forecast System (IFS) of the European Centre for Medium-Range Weather Forecasts (ECMWF).

NWP has continuously improved over the seven decades since the first barotropic model forecast (Benjamin et al., 2019). Current state-of-the-art models typically provide skillful predictions of global weather patterns at effective grid point spacings of roughly 0.1° of latitude (about 10 km) through at least seven days of forecast lead time (Bauer et al., 2015). The computational effort required to generate such global high-resolution

68 forecasts is enormous and only available at a handful of advanced dedicated centers. En-
 69 semble forecasts, which provide an important way to account for uncertainty by gener-
 70 ating a set of equally plausible predictions and extend the limit of skillful forecasts be-
 71 yond that of a single deterministic model run, are also limited by the computational bur-
 72 den of high-resolution NWP to about 50 members (Palmer, 2019).

73 Global NWP models represent 3D fields as sets of nested spherical shells in which
 74 the distance between each shell is the local vertical grid spacing. On every time step, the
 75 ECMWF Integrated Forecasting System (IFS), as configured for sub-seasonal forecast-
 76 ing, updates 10 prognostic 3D variables defined at 91 vertical levels. Along with surface
 77 pressure, this totals to over 900 spherical shells of data. Here, we use “spherical shell of
 78 data” to describe a single variable defined at a single vertical level on a spherical shell
 79 covering the globe. The large number of spherical shells of data (combined with the fine
 80 horizontal resolution) in NWP models is required to produce acceptably accurate numer-
 81 ical solutions to the equations governing atmospheric motions. The data at each indi-
 82 vidual point, however, cannot be independently perturbed while maintaining a meteo-
 83 rologically relevant atmospheric state. For example, on horizontal scales larger than about
 84 10 km, the temperatures throughout a vertical column and the heights of constant pres-
 85 sure surfaces must satisfy hydrostatic balance.

86 The actual number of independent degrees of freedom required to represent the pre-
 87 dictable components of the global atmosphere is unknown, but it clearly decreases with
 88 increasing forecast lead times (Lorenz, 1969). GraphCast (Lam et al., 2022), for exam-
 89 ple, has achieved success at lead times as short as 6 h with 227 spherical shells of data.
 90 It can produce forecasts using much less computation time than the ECMWF IFS, but
 91 it still requires large computing resources for training: 3 weeks using 32 TPU 4 proces-
 92 sors. Pangu-Weather (Bi et al., 2023) cuts the number of spherical shells by almost 2/3
 93 to 69. The spherical Fourier neural operator (SFNO) version of FourCastNet compared
 94 with the IFS in (Bonev et al., 2023) uses 73 spherical shells of data. Here, we take this
 95 reduction much farther, presenting a parsimonious DLWP model that uses just 7 spher-
 96 ical shells of data to efficiently provide forecasts approaching the skill of ECMWF. While
 97 not as accurate as GraphCast or Pangu-Weather for medium range forecasts with lead
 98 times less than two weeks, we demonstrate that our model generates far less bias in fore-
 99 casts of 500 hPa height in one-year iterative forecasts. In addition, our model is poten-
 100 tially better suited for research applications such as computing the sensitivities of its com-
 101 pact state vector to custom diagnostic functions by backpropagation.

102 In contrast to many of the recent DLWP architectures, our approach relies on con-
 103 volutional neural networks (CNN), building on early work by Scher and Messori (2018)
 104 and Weyn et al. (2019) and the U-Net configuration in Weyn et al. (2020) and Weyn et
 105 al. (2021). Here, we document substantial improvements over Weyn et al. (2021), obtained
 106 by replacing the cubed sphere data representation with the HEALPix mesh, which is widely
 107 employed in astronomy (Gorski et al., 2005). In addition, we improve the former model
 108 by implementing physically motivated modifications in form of residual connections, re-
 109 current modules, and inverting the channel depth as compared with a standard U-Net.

110 2 Related Work

111 Pioneering efforts to create machine learning models to forecast the weather from
 112 reanalysis or general circulation model (GCM) output include the dense neural network
 113 of Dueben and Bauer (2018) and the CNN models of Scher and Messori (2019) and Weyn
 114 et al. (2019), all of which employed latitude longitude (lat-lon) meshes. Weyn et al. (2020)
 115 obtained significantly improved forecasts by switching to a cubed sphere mesh with a
 116 CNN in the standard U-Net architecture (Ronneberger et al., 2015). Their model was
 117 capable of generating realistic weather patterns when stepped forward for a full year (730
 118 12 h steps). Retaining the cubed sphere, Weyn et al. (2021) produced forecasts out to

119 sub-seasonal time scales using large multi-model ensembles, and Lopez-Gomez et al. (2022)
 120 migrated from the U-Net into a U-Net 3+ architecture (Huang et al., 2020)—which adds
 121 connections between multiple hierarchical levels in the U-Net—to generate forecasts of
 122 extreme surface temperatures.

123 Returning to the lat-lon mesh, Rasp and Thuerey (2021) demonstrated that a deep
 124 Resnet could be pre-trained on GCM data and then fine-tuned by transfer learning on
 125 ERA5 data to produce up to 5-day forecasts at coarse 5.65° grid spacing. Building on
 126 transformer models from computer vision (Dosovitskiy et al., 2020; Guibas et al., 2021),
 127 Pathak et al. (2022) and Kurth et al. (2022) used Fourier neural operators (Li et al., 2020)
 128 to develop FourCastNet on a 0.25° lat-lon mesh to generate forecasts approaching the
 129 accuracy of ECMWF’s IFS. FourCastNet was not, however, capable of stable long-lead-
 130 time autoregressive rollouts. This difficulty was overcome by switching from 2D Fourier
 131 modes on a lat-lon mesh to spherical harmonic functions Bonev et al. (2023). The result-
 132 ing SFNO model eliminated much of the vision transformer architecture while improv-
 133 ing accuracy and remaining stable for one-year forecasts.

134 Again on a 5.65° lat-lon mesh, Hu et al. (2022) used a shifted window (Swin) trans-
 135 former (Liu et al., 2021) to produce single forecasts as well as ensembles generated by
 136 perturbing the latent state using samples from their learned distribution. Bi et al. (2023)
 137 also applied Swin transformers on a lat-lon mesh, but used a fine 0.25° lat-lon grid spac-
 138 ing, 3D transformers, and included latitude and longitude fields as input to train a “3D
 139 Earth-specific transformer” at four different forecast lead times of 1, 3, 6, and 24 h, which
 140 are used in combination to span an arbitrary hourly forecast period with minimal model
 141 steps. If the ECMWF IFS NWP forecasts are averaged to the coarser 0.25° lat-lon mesh,
 142 Pangu-Weather outperforms NWP on several metrics.

143 In contrast to the preceding approaches, graph neural networks (Gori et al., 2005;
 144 Scarselli et al., 2008; Kipf & Welling, 2016; Battaglia et al., 2018; Pfaff et al., 2020) were
 145 applied on icosahedral meshes at course resolution by Keisler (2022) and at fine resolu-
 146 tion in the Graphcast model (Lam et al., 2022). Similarly to Pangu-Weather, GraphCast
 147 appears to outperform the coarsened ECMWF IFS forecast on several metrics.

148 3 Methods

149 3.1 Data

150 3.1.1 Choice of Variables

151 Beginning with the same six prognostic variables used in Weyn et al. (2021)—geopotential
 152 height at 1000 hPa and 500 hPa (Z_{1000} , Z_{500}),¹ 700 hPa to 300 hPa thickness ($\tau_{700-300}$)
 153 defined as $Z_{300} - Z_{700}$, temperature at 2 m height above ground (T_{2m}), temperature at
 154 850 hPa (T_{850}), and total column water vapor ($TCWV$)—we add Z_{250} based on its im-
 155 portance in the model of Rasp and Thuerey (2021) and to provide an upper tropospheric
 156 variable. As in Weyn et al. (2021), three prescribed fields are also provided: topographic
 157 height, land-sea mask, and top-of-atmosphere (TOA) incident solar radiation. We do not
 158 include prescribed or predicted sea-surface temperature or surface fluxes above the land
 159 or ocean. *No* specific information about position on the globe, such as latitude and lon-
 160 gitude, is provided. Three-hourly data from the ERA5 reanalysis (Hersbach et al., 2020)
 161 provide training data from 1979-2012, a validation set from 2013-2016, and a test set from
 162 2017-2018.

¹The related variable in the ERA5 dataset is geopotential and named z , whereas the geopotential height, typically referred to as Z , represents the actual height above sea level of the respective pressure surface and is obtained by dividing geopotential by the gravitational constant.

163 **3.1.2 HEALPix Mesh**

164 We discretize all fields using the Hierarchical Equal Area isoLatitude Pixelization
 165 (HEALPix) (Gorski et al., 2005). As depicted in Figure 1, a HEALPix mesh is formed
 166 by dividing the sphere into twelve equal-area diamond-shaped faces, with four faces ly-
 167 ing in the northern and southern hemispheres, and four in the tropics. According to Gorski
 168 et al. (2005), the HEALPix mesh has three important properties. (1) *Hierarchical struc-*
 169 *ture of the database:* Each of the twelve base faces can be progressively subdivided into
 170 smaller patches. (2) *Equal areas for the discrete elements of the partition:* All patches
 171 are the same size. (3) *Isolatitude distribution for the discrete area elements on the sphere:*
 172 The patches line up with lines of latitudes, facilitating the computation of zonal averages
 173 and one-dimensional zonal spectra. Importantly, this last property makes the HEALPix
 174 mesh an “east is to the right” grid, which facilitates the training of CNN kernels to cap-
 175 ture the motion of typical weather disturbances, as discussed in subsection 4.1.

176 The HEALPix can be considered a graph and does not allow a seamless applica-
 177 tion of convolution operations. Thus, Perraudin et al. (2019) explicitly define a graph
 178 from the HEALPix—by connecting adjacent neighbors with weighted edges—and per-
 179 form a graph convolution to classify weak lensing maps from cosmology. In a different
 180 approach, Krachmalnicoff and Tomasi (2019) classify digits and determine cosmic param-
 181 eters from simulated cosmic microwave background maps. They apply 1D convolutions
 182 to the flattened HEALPix data with a kernel size k and stride s both equal to 9, append-
 183 ing a zero to those cases where only seven instead of eight neighbors are defined (top cor-
 184 ner of the tropical faces). In contrast, we treat the twelve faces as distinct images and
 185 pad their boundaries using data from neighboring faces to allow the computation of 2D
 186 convolutions and averaging operators directly, as detailed in section Appendix A. To ac-
 187 celerate the padding operation, we have implemented a custom CUDA kernel, which is
 188 available in our repository.²

189 The grid spacing, or shortest inter-node spacing, on the HEALPix mesh is the di-
 190 agonal distance between a pair of nodes on adjacent latitude lines. Denoting a HEALPix
 191 mesh with n divisions along one side of the original 12 faces as HPX n . The grid spac-
 192 ing is approximately 220 km ($\approx 2^\circ$) for HPX32 and 110 km ($\approx 1^\circ$) for HPX64.³

193 **3.2 Machine Learning Architecture**

194 Relating to Tobler’s first law of geography: “All things are related, but nearby things
 195 are more related than distant things.” (Tobler, 1970), we mostly retain the comparably
 196 simple U-Net structure from Weyn et al. (2020). U-Nets (Ronneberger et al., 2015) are
 197 hierarchically structured feed-forward convolutional neural networks that were originally
 198 proposed for segmenting biomedical images. The U-Net structure proposed here intro-
 199 duces several physically motivated advancements to the vanilla U-Net used by Weyn et
 200 al. (2021) for time-series forecasting. The advancements and model configurations are
 201 visualized in Figure 2, detailed in Table B1, and described in the following.

202 **3.2.1 Residual Prediction**

203 We switch to a residual prediction approach both for the full predictive step and
 204 within each ConvNeXt block.⁴ Predicting changes over a time step, instead of the full
 205 fields, is similar to the discretization of time derivatives when solving partial or ordinary

² <https://github.com/CognitiveModeling/dlwp-hpx>

³ We provide download explanations and projection scripts in our repository. The 3D HEALPix figures are drawn in Blender 3.4.1; respective Blender files are provided in the repository too.

⁴ As detailed in Figure 2, we modify the original ConvNeXt block from Liu et al. (2022) by removing the bottleneck and employing a two-stage convolution as done in Weyn et al. (2021).

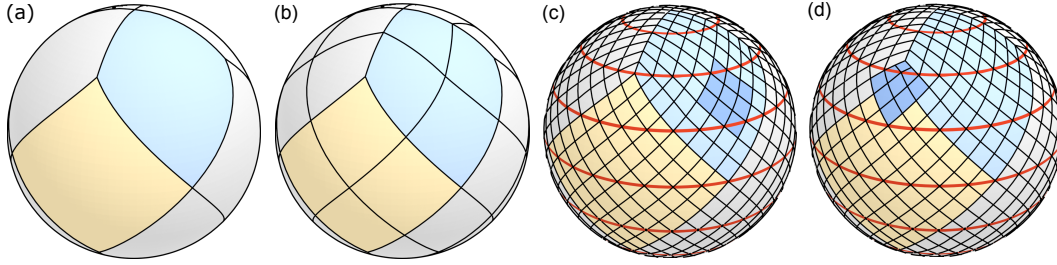


Figure 1: Division of the sphere into twelve faces according to the HEALPix. Four faces to represent either the northern (blue) and southern extratropics, while four more faces arrange around the equator to represent the tropics (yellow). Each face can be subdivided into patches with divisions along the side of each face given by powers of two. The sphere in (a) has a pixel-count of one per face side; we call it **hpx1**. The sphere in (b) counts two pixels per side (**hpx2**), whereas the two spheres in (c) and (d) have eight pixels per side, i.e., **hpx8**. Several latitude lines in red emphasize the iso-latitudinal arrangement of the patches. The saturated blue area depicts a 3×3 stencil, as applied by a standard convolution. To apply the 3×3 stencil at the top corner of the equatorial faces, i.e., stencil position in (d), we simulate a hypothetical patch by computing the average from the according extratropical face patches.

206 differential equations, and has been used successfully in previous deep-learning weather
 207 prediction models (Pathak et al., 2022; Keisler, 2022; Hu et al., 2022; Lam et al., 2022).

208 **3.2.2 Inverting the Ordering of Channel Depth**

209 The standard U-Net for semantic segmentation (Ronneberger et al., 2015) and its
 210 successors (Zhou et al., 2018; Huang et al., 2020) employ relatively few channels on the
 211 highest level and successively double the channel depth, while halving the spatial reso-
 212 lution in each deeper layer. This ordering is useful in image segmentation tasks, where
 213 deeper channels are required to create increasingly abstract filters to identify semantic
 214 features and express complex objects. In weather prediction, however, we find it is bet-
 215 ter to devote more capacity to the layers in the first level, where a wide variety of fine
 216 grained weather phenomena must be captured. Deeper layers at coarser resolution, on
 217 the other hand, need only encode larger scale atmospheric motions, which can be ade-
 218 quately represented with comparably fewer channels.

219 Thus, we invert the channel order, employing 136, 68, and 34 channels in each con-
 220 volution on the first, second, and third layer, respectively (cf. Figure 2). While this mod-
 221 ification improves the model performance significantly, it also increases the computational
 222 burden, since more computations and data processing are required to evaluate the ad-
 223 ditional convolutions at fine spatial resolution. Tests which preserved the total number
 224 of trainable parameters, but completely eliminated the deeper layers in the U-Net gave
 225 worse results, demonstrating that the longer-range connections and richer latent space
 226 structures enabled by the full U-Net architecture remain important.

227 **3.2.3 Recurrent Modules**

228 The vanilla U-Net is a feed-forward network, which treats successive inputs inde-
 229 pendently even if the data represents a continuous sequence over time. Feed-forward net-
 230 works do not have any memory capacity. They do not maintain an internal state between
 231 time steps. To enable the exploitation of information from previous latent states, we in-
 232 clude a gated recurrent unit (GRU) (Cho et al., 2014) at the end of each decoder block

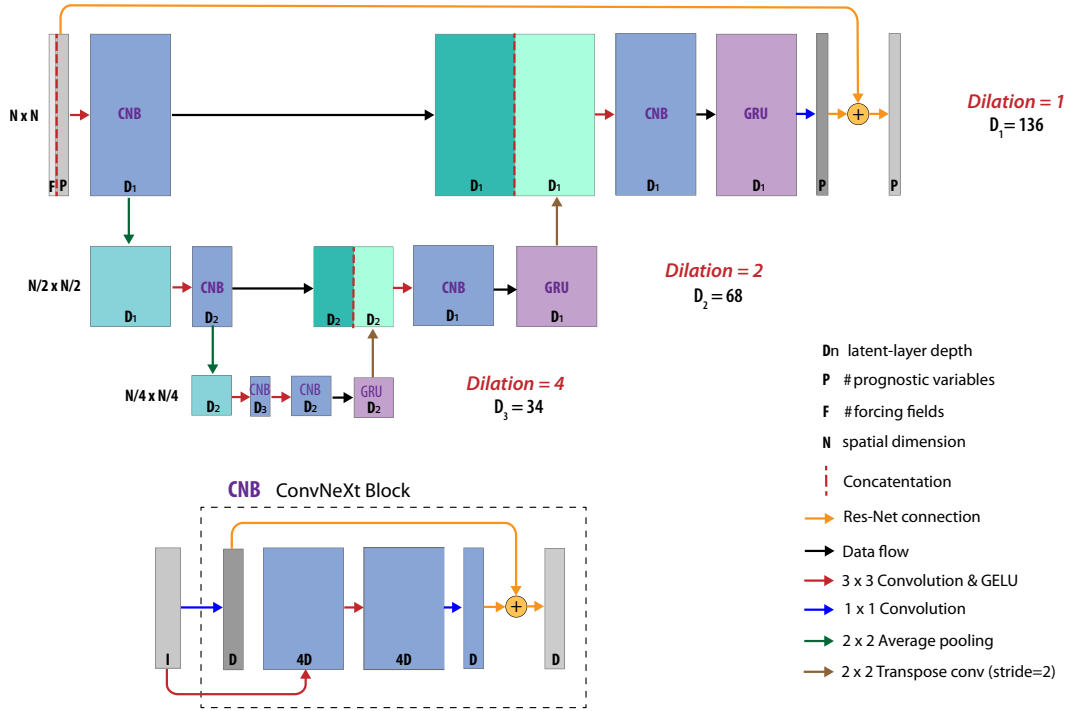


Figure 2: Schematic representation of the DLWP-HPX architecture for our best performing model. There is one ConvNeXt block at each level in both the encoder and the decoder. In contrast to the configuration in typical image processing applications, the channel, or latent-layer, depth decreases from 136 to 68 to 34 at deeper layers in the U-Net.

233 with kernel size $k = 1$. We chose GRUs over LSTMs (Hochreiter & Schmidhuber, 1997)
 234 since we re-initialize the recurrent data over each 24 h-cycle, and therefore do not require
 235 forget-gates (as confirmed experimentally, not shown).

236 3.2.4 Miscellaneous Modifications

237 Several other components of the original Weyn et al. (2021) model were modified
 238 based on recent results from deep learning research: the capped leaky ReLU was replaced
 239 by capped GELU activations (Hendrycks & Gimpel, 2016); upsampling was changed from
 240 nearest-neighbor sampling (knn-sampling with $k = 1$) to a transposed convolution; fi-
 241 nally, the pairs of two successive convolutions were replaced at each encoder and decoder
 242 level in the U-Net by a modified ConvNeXt block (Liu et al., 2022), as visualized in Fig-
 243 ure 2.

244 3.2.5 Time Stepping Scheme

245 Similarly to Weyn et al. (2021), we apply a two-in-two-out mapping with a tempo-
 246 ral resolution twice as fine as the actual time step. For example, two atmospheric states
 247 3 h apart (each consisting of seven prognostic, along with three prescribed fields) are con-
 248 catenated and input to the model, which generates a new pair of states, each character-
 249 ising the atmosphere 6 h later in time. This strategy is observed to stabilize and accel-
 250 erate the training, since the model receives additional information about the atmosphere’s
 251 rate of change and only has to be called half as often.

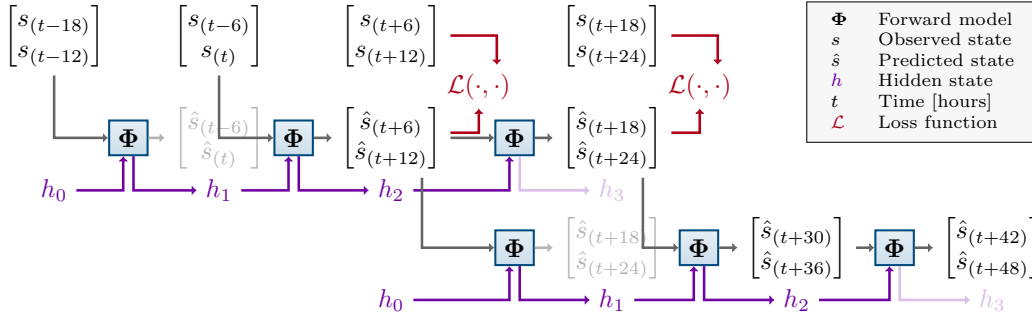


Figure 3: Two time-level input-output scheme with GRU for training and inference assuming 6 h time resolution. The output from the preliminary initialization step (light gray) is discarded, but the hidden state h_1 is generated and used in the first model step. The hidden state h_3 (light purple) at the end of the 24 h forecast is discarded as the GRU will be re-initialized for the next recursive inference step (lowest row). For training (top right), the loss function is computed from the four forecast times spanning 24 h period at 6 h resolution, as indicated in red.

252 The frequency spectrum of atmospheric kinetic energy has a strong peak at 24 h
 253 because many circulations are modulated by solar heating. We therefore evaluate the
 254 training loss function as the mean squared error over a 24 h period. Tests in which the
 255 MSE was evaluated over multi-day periods tended to result in a model that gradually
 256 approached climatology over many recursive steps.

257 Training our model only over one daily cycle does mean that the recurrent states
 258 of the GRUs are not optimized for long rollouts. To prevent the explosion of recurrent
 259 states when generating long multi-day forecasts, we re-initialize the recurrent states ev-
 260 ery 24 h as illustrated in Figure 3 for a 12 h time step with 6 h resolution. For training
 261 or for the first step in a long forecast rollout, the model predicts $[\hat{s}_{(t+6)}, \hat{s}_{(t+12)}]$ from ini-
 262 tial data $[s_{(t-6)}, s_{(t)}]$, and then in the subsequent step uses $[\hat{s}_{(t+6)}, \hat{s}_{(t+12)}]$ to predict $[\hat{s}_{(t+18)}, \hat{s}_{(t+24)}]$.
 263 But before this, the hidden states for the GRUs are initialized in a preliminary step by
 264 calling the model once with the state pair $[s_{(t-18)}, s_{(t-12)}]$ and a hidden state h_0 initial-
 265 ized with zeros. The resulting forecast for $[\hat{s}_{(t-6)}, \hat{s}_{(t)}]$ is discarded, but the hidden state
 266 h_1 is supplied to the GRU and paired with the actual initial data $[s_{(t-6)}, s_{(t)}]$ for the first
 267 step of the model. As shown by the bottom row in Figure 3, in a forecast rollout, the
 268 next day’s prediction begins by re-initializing the GRU starting with forecast values from
 269 one time step earlier and h_0 set to zero to obtain h_1 . Note that since the GRU is re-initialized
 270 every day, there would be five model steps per day when using a 6 h time step (with 3 h
 271 data resolution).

272 3.2.6 Training

273 Our best performing DLWP-HPX model, described above, has 9.8 M parameters
 274 that are trained for 300 epochs (equivalent to 931,199 update steps) over eight days on
 275 four NVIDIA A100 GPUs with 80 GB VRAM each. A batch size of eight per GPU is
 276 chosen, effectively resulting in an overall batch size of 32. We combine the Adam opti-
 277 mizer (Kingma & Ba, 2014) with a cosine annealing learning rate scheduler (Loshchilov
 278 & Hutter, 2016), setting the initial learning rate to 2×10^{-4} and gradually refining it
 279 to zero. To stabilize the training, we clip the gradients to the current learning rate, which
 280 we observe to be particularly beneficial for large and recurrent models.

3.3 The Receptive Field

Several leading DLWP models (Pathak et al., 2022; Hu et al., 2022; Bi et al., 2023; Chen et al., 2023) are based on Vision Transformers (ViTs) (Dosovitskiy et al., 2020), which were originally developed to account for non-local relationships in images; effectively working on patch embeddings. ViTs are successors of Transformers (Vaswani et al., 2017), which were introduced to efficiently accommodate very non-local relationships in natural language processing (NLP), where no fixed upper bound exists on the distance between words that may interact to change the meaning of a sentence. In contrast to ViTs, we use a U-Net to emphasize local atmospheric interactions, nevertheless each step of our model samples from a very large receptive field.

There is a strong physical constraint on the locality of atmospheric interactions, which is that *no atmospheric disturbances travel faster than the speed of sound*, roughly 300 m/s. Sound waves are not meteorologically significant, and are not represented in the data used to train ML weather models. A better measure of the speed of the fastest moving signals of meteorological importance is the transport by the strongest jet-stream winds, which could transport a passive tracer at roughly 100 m/s, or about 4300 km in 12 h.

The pair of 2×2 average poolings and the dilations in the second and third levels of our U-Net architecture (Figure 2) substantially widen the receptive field that potentially influences the solution at a given point after each forward step of our model. Neglecting influences from special points at the corners of the twelve basic HEALPix faces, the receptive field at each stage of the neural network is listed in Table B1 and grows to a 175×175 patch of cells after the last 3×3 convolution in the decoder.

The diagonal distance between adjacent points on our 3×3 stencil (dark blue patch in Figure 1) on a HPX64 mesh is approximately 110 km. Thus, the receptive field for one step of our full HPX64 model is a patch exceeding 18 900 km on each side, which is large enough to include all points influenced by sound wave propagation over a 12 h time step, and far more than would be required to contain the fastest moving meteorologically significant signals present in the ERA5 training data. In particular, at every step, our HPX64 forecast at a given point is influenced by a set of surrounding points containing roughly 70% of all the cells covering the globe.

4 Results

In the following, we first evaluate several key variables in our model over a 14-day forecast lead time, which is slightly longer than the period over which knowledge of the initial atmospheric conditions gives these single deterministic forecasts some predictive skill. We compare our best model with the ECMWF S2S forecasts and with our previous Weyn et al. (2021) results. We then document the successive improvements that our changes in model architecture have on the RMSE and ACC scores for Z_{500} . Next, we examine the ability of the model to distinguish between the amplitudes of the daily T_{2m} ranges in tropical forests, in deserts, and over the ocean. Finally, we examine the behavior of the simulations over sub-seasonal (eight-week) and one-year free running rollouts.

4.1 Quantitative Performance Through 14-Day Forecast Lead Time

To compare our model with the results from Weyn et al. (2021) and to state-of-the-art NWP from ECWFMF, we compute both root mean squared error (RMSE) between observations and model predictions and anomaly correlation coefficient (ACC) scores with respect to the ERA5 climatology. Both metrics are compared on a $1^\circ \times 1^\circ$ lat-lon mesh and weighted by latitude, requiring us to project our DLWP-HPX and Weyn et al. (2021) forecasts from the HEALPix and cube-sphere meshes onto the lat-lon grid. Because our

Table 1: Number of trainable parameters in millions, number of spherical shells of prognostic variables, horizontal resolution in degrees latitude, and temporal resolution (Δ_t) of the models compared in Figure 4.

Model	Parameters	Spherical shells	Resolution	Δ_t
Weyn 2021	2.7M	6	1.4°	6 h
Our HPX64	9.8M	7	1°	3 h
ECMWF	—	900+	0.15°	0.2 h

329 ultimate focus is on sub-seasonal and seasonal forecasting, we compare against ECMWF’s
 330 integrated forecasting system for sub-seasonal forecasts (IFS S2S), which were initialized
 331 bi-weekly on Mondays and Thursdays and stepped forward at about 16 km effective res-
 332 olution for the first 15 days (then doubling to 32 km).⁵ For comparison with Weyn et
 333 al. (2021), our test set focuses on the years 2017 and 2018. In this and all the following
 334 cases, except a few simulations in our ablation study, computations are performed at HPX64
 335 and 3 h resolution (corresponding to 6 h time steps). Key parameter attributes of the model
 336 from Weyn et al. (2021), IFS S2S, and our HPX64 model are listed in Table 1.

337 As shown in Figure 4, the RMSE scores for Z_{500} , 24-hour-averaged T_{2m} (instan-
 338 taneous T_{2m} fields are not archived from the ECMWF S2S forecasts⁶) and T_{850} all im-
 339 prove substantially compared to Weyn et al. (2021). Moreover, despite the small num-
 340 ber of prognostic variables and coarse spatial resolution of our model, the RMSEs for
 341 Z_{500} and T_{850} only lag the scores for ECMWF S2S by about 1 day at one-week lead time.
 342 As expected theoretically, the RMSE scores for all three models appear to be asymptot-
 343 ically approaching $\sqrt{2}$ times climatology beyond two weeks when the skill of a single de-
 344 terministic forecast drops toward zero. We present the comparison of 24-hour-averaged
 345 T_{2m} between our model and IFS S2S for completeness, but it should be interpreted with
 346 caution. The re-gridding of both the IFS S2S and the HEALPix data to the $1^\circ \times 1^\circ$
 347 lat-lon analysis grid introduces errors in the representation of coastlines and topography
 348 that significantly influence the surface temperature field. As a consequence, the RMSE
 349 values shown in Figure 4 (b) are not representative of those in each model’s native rep-
 350 resentation of the T_{2m} field.

351 One additional issue that arises when plotting initial RMSE (and to a lesser extent
 352 ACC) for the ECMWF IFS S2S model is that, unlike our DLWP-HPX model, the IFS
 353 forecasts are not initialized with the ERA5 data. Thus, at very short forecast lead times,
 354 differences between the IFS initialization and the ERA5 data introduce apparent errors
 355 in the IFS forecast that are not representative of its actual performance. Lam et al. (2022)
 356 accounted for this in their comparison between the IFS and GraphCast, but it requires
 357 considerable extra computation. We are not claiming to outperform the IFS, so we sim-
 358 ply suggest using caution when comparing errors between our models and the IFS at lead
 359 times less than 2 days.

360 ACC scores for Z_{500} , T_{2m} , and T_{850} are also shown in Figure 4(d)–(f). As with RMSE,
 361 there is substantial improvement relative to both the previous model from Weyn et al.
 362 (2021) and the IFS S2S. In meteorological contexts, an ACC score of 0.6 is typically con-
 363 sidered the lower limit of practical skill. The scores from our HEALPix model cross this

⁵ <https://confluence.ecmwf.int/display/S2S/ECMWF+model+description>

⁶ <https://apps.ecmwf.int/datasets/data/s2s-realtime-daily-averaged-ecmf/levtype=sfc/type=cf/>

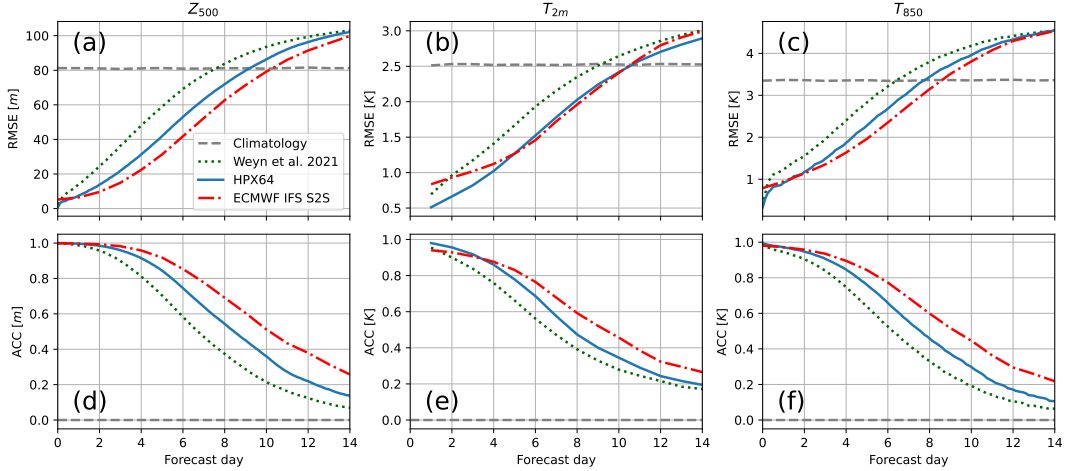


Figure 4: Comparison of the performance of the DLWP-HPX, Weyn et al. (2021), and ECMWF IFS S2S models, averaged over 208 forecasts. RMSE for (a) Z_{500} , (b) T_{2m} and (c) T_{850} ; climatology is indicated by the gray dashed line. ACC for (d) Z_{500} , (e) T_{2m} and (f) T_{850} .

364 threshold at about 7.5 days for Z_{500} and 6.5 days for T_{850} , both of which are about 1.5
 365 day sooner than the respective results for the IFS S2S. Numerical comparisons of the model
 366 RMSE and ACC scores averaged over the same 208 forecasts used to plot Figure 4 are
 367 given for 3-day and 5-day lead times in Table 2.

368 The relative importance of the various improvements in model architecture between
 369 Weyn et al. (2021) and our best DLWP-HPX model is illustrated for the Z_{500} field in
 370 Figure 5. The total number of trainable parameters is held constant at roughly $2.7 \times$
 371 10^6 over the first five sets of changes. The RMSE rises to 50 m around 4.2 days in Weyn
 372 et al. (2021) (dark green dotted curve); replacing the 64×64 cubed sphere by a HPX32
 373 grid (aqua curve) delays the error growth by about 0.5 day despite the associated 50%
 374 reduction in total grid points. There is also a similar substantial improvement in the ACC.
 375 Continuing with the HPX32 mesh, we replace the capped ReLU by a capped GELU acti-
 376 vation function, replace knn-interpolation by strided transposed convolution, and in-
 377 troduce dilated convolutions in the two lower levels of the U-Net (as detailed in Figure 2);
 378 this yields the modest but distinct improvements shown by the dark-blue curves.

379 Next, we replace the pairs of convolutions in each level of the encoder and decoder
 380 by a ConvNeXt block without a bottleneck (dashed tan curve). This actually produces
 381 a slight degradation in performance, but in other configurations closer to our final model,
 382 the ConvNeXt block does improve the performance, and importantly, it also reduces the
 383 memory footprint by about 25% at a constant parameter count. A further significant im-
 384 provement is obtained by inverting the standard U-Net progression in channel depth to
 385 have the most channels at the highest spatial resolution and the fewest at the lowest re-
 386 solution (dark red curve). The final significant improvement in the 2.7-million param-
 387 eter model is obtained by adding recurrence in the form of GRU cells in the decoder (green
 388 curve).

389 After adding the GRU cells, the rise of the RMSE to 50 m is delayed to about 5.3
 390 days and the drop of the ACC below 0.6 to roughly 6.8 days. The next series of changes
 391 produces successive small improvements that push these values out to about 5.7 days
 392 for RMSE and 7.4 days for ACC. These improvements, as sequentially plotted in Fig-
 393 ure 5, are: increasing the number of trainable parameters to 9.8×10^6 , adding the Z_{250}

Table 2: Root mean squared errors (RMSE) and anomaly correlation coefficient (ACC) scores for Weyn et al. (2021) (W21), our HPX64, and ECMWF’s IFS models, evaluated on geopotential at 500 hPa (Z_{500}), temperature 2 m above ground (T_{2m}), and temperature at 850 hPa (T_{850}) on lead times of 3 and 5 days.

	Lead time	Z_{500}			T_{2m}			T_{850}		
		W21	HPX64	IFS	W21	HPX64	IFS	W21	HPX64	IFS
RMSE	3 days	36.26	21.88	14.91	1.17	0.82	1.02	1.95	1.49	1.35
	5 days	59.01	41.91	31.30	1.67	1.27	1.27	2.83	2.28	1.96
ACC	3 days	0.90	0.96	0.98	0.84	0.92	0.91	0.84	0.91	0.94
	5 days	0.70	0.84	0.92	0.66	0.78	0.83	0.64	0.76	0.84

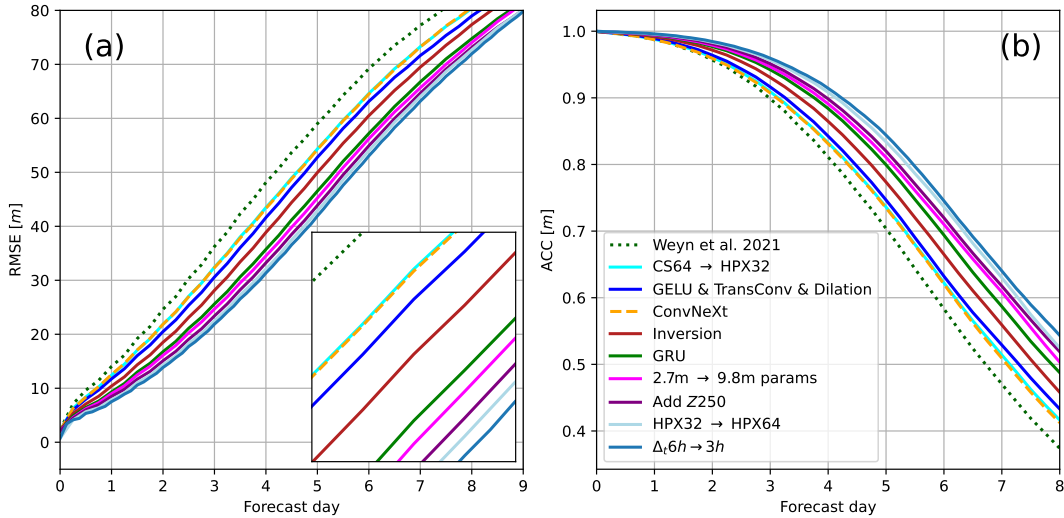


Figure 5: Impact of successive model improvements on the accuracy of Z_{500} building from WDC to our HPX64 model with $\Delta_t = 3h$. Each successive change builds on top of the previous architecture, adding the modification indicated in the legend: (a) RMSE, (b) ACC. Inset in (a) provides a magnified view of the error growth between 5 and 6 forecast days.

394 field, increasing the horizontal resolution to HPX64 (which is more important for ACC
 395 than RMSE particularly on T_{2m}), and decreasing the time resolution to 3h. Benefits from
 396 the use of 3h time resolution were only obtained if the model was configured with the
 397 GRUs.

398 The single most effective modification in the preceding set of successive improve-
 399 ments is the migration from the cubed sphere to the HEALPix mesh, even though the
 400 64×64 cubed sphere has twice the total number of grid-points as the HPX32 mesh.
 401 The most likely explanation for the superiority of the HEALPix mesh is not simply that
 402 it is a more uniform covering of the globe than that provided by the cube-sphere, but
 403 that east and west have the same orientation in every HEALPix cell; we refer to this prop-
 404 erty as “east to the right.” In particular, the center and the east and west corners of each
 405 HEALPix cell are all at the same latitude. (A similar relationship holds in the north-south
 406 direction for meridians passing through those cells lying equatorward of the maximum
 407 north-south extent of the four equatorial faces in Figure 1 (a).) Thus, on the HEALPix
 408 mesh, eastward motion at all points and at all latitudes would be in the same direction

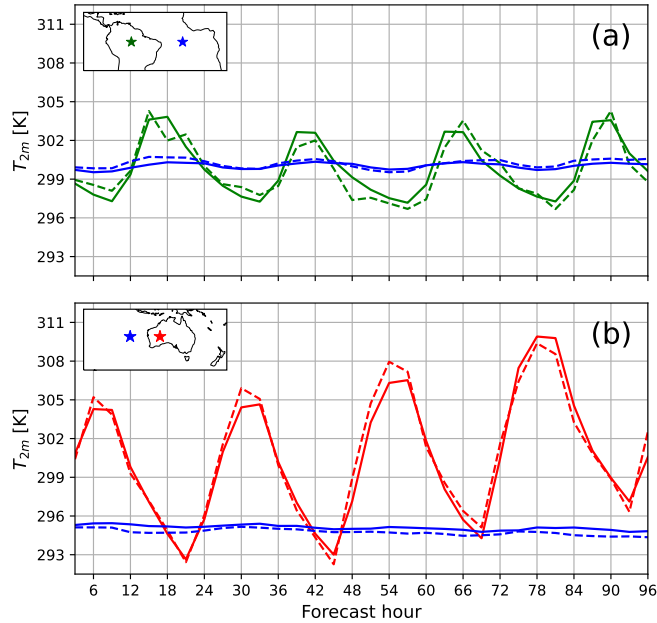


Figure 6: HPX64 simulation of the diurnal cycle of T_{2m} (solid curves) at the four locations shown in the insets starting from 00 UTC on 12 March 2018. ERA5 values for the same $1^\circ \times 1^\circ$ lat-lon cell are shown as dashed lines. Values are plotted every 3 h.

409 across the diamond-shaped 3×3 stencil in Figure 1 (c). In contrast, at any point on either
 410 of the polar faces on the cubed sphere, east could map to any of four directions along
 411 the axes of the 3×3 convolutional stencil, depending on its longitude, as visualized in
 412 section Appendix A.

413 In mid- and high-latitudes, most large-scale weather systems move in a generally
 414 eastward direction. We believe this allows a fixed number of kernel elements to more ef-
 415 ficiently and effectively produce the required set of flow evolutions in the latent layers.
 416 To a lesser extent, this same consideration also applies to the four equatorial faces of the
 417 cubed sphere, where, for example, eastward flow near the northeastern corner of a face
 418 would need to move at an angle relative to the northern side of the stencil that is oppo-
 419 site in sign to that required in the northwestern corner.

420 4.2 Eliminating the Need for Boundary-Layer Parameterizations

421 Accurate forecasts of surface temperatures in NWP models rely on the empirical
 422 parameterization of multi-scale processes near the Earth’s surface in the atmospheric bound-
 423 ary layer (ABL). The bottom of the ABL includes the roughness layer (2–5 times the
 424 height of roughness elements such as vegetation), and the surface layer (often 10–100 m
 425 deep), where shear-driven turbulence dominates generation by convection. The depth
 426 of the full ABL, where larger-scale eddies and circulations communicate the processes
 427 in the surface layer to the free atmosphere, can vary from $O(100)$ m in calm stable night-
 428 time conditions to several kilometers during the day over deserts.

429 No effort is made to explicitly account for ABL processes in our model; the T_{2m}
 430 field is treated the same as the other six prognostic fields. The same CNN kernels are
 431 employed everywhere over the globe on the HEALPix mesh; the only data that might
 432 distinguish one location from another are the land-sea mask, the terrain elevation, and
 433 the TOA solar forcing; neither longitude nor latitude are provided. Yet our model does
 434 a good job of capturing the diurnal cycle in multi-day forecasts over very different sur-
 435 faces. Figure 6 shows the diurnal cycle in T_{2m} at locations over the Amazon forest, the
 436 Australian desert, and two adjacent oceans over a 4-day simulation starting at 00 UTC
 437 on 12 March 2018.

438 Compared to over land, the diurnal T_{2m} variations are modest over the oceans, and
 439 they are well captured by our model. The land-sea mask is undoubtedly important in
 440 distinguishing the ocean locations from those over land. More interestingly, the model
 441 does an excellent job of capturing the large diurnal temperature range over the Australian
 442 desert, while correctly generating a much lower amplitude signal over the Amazon. The
 443 prognostic field that has most likely facilitated this distinction is $TCWV$, which is sig-
 444 nificantly higher over the Amazon than over the Australian desert. The model also cap-
 445 tures the 4-day trend for increasing temperatures over Australia, which is linked to the
 446 evolution of larger-scale weather systems. Overall, the ability of the model to capture
 447 the diurnal T_{2m} cycle with just seven prognostic fields, without any special treatment
 448 of the ABL, and without geo-specific inputs such as latitude and longitude is suggestive
 449 of the power and potential of DLWP-HPX.

450 4.3 Iterative Rollouts Over Subseasonal to Annual Time Scales

451 There are three time scales of primary interest for global atmospheric simulations:
 452 medium-range weather forecasting for lead times of up to two weeks, sub-seasonal and
 453 seasonal forecasts for lead times up to 6–9 months, and climate simulations over periods
 454 of tens to hundreds of years. Our focus is on the sub-seasonal to seasonal time scale; there-
 455 fore, in this section we examine the model’s performance in iterative rollouts over peri-
 456 ods up to one year.

457 To investigate the stability and drift in model simulations over a full annual cycle,
 458 we initialize it using ERA5 data for 00 UTC on 1 June 2017 (together with the 21 UTC
 459 fields on 31 May). Using 6 h time steps (with 3 h time resolution), we perform 1460 it-
 460 erations to generate a 365-day simulation. The three-day running mean of Z_{500} , aver-
 461 aged around each latitude, is plotted as a function of latitude and time in Figure 7, along
 462 with the corresponding averages from the ERA5 data. Despite being trained to minimize
 463 RMSE over a single day and not enforcing any physical constraints, the DLWP-HPX sim-
 464 ulation responds to the TOA solar forcing to generate the annual cycle reasonably well.

465 One region where the errors are significant is the arctic. About 5 months into the
 466 simulation, the simulated heights in the arctic region drop as much as 60 m below those
 467 in the reanalysis during the boreal winter. In contrast, at 5–8-month lead times, the heights
 468 in the antarctic region increase to approximately correct values in the austral summer.
 469 The asymmetry between the response in arctic and antarctic flips if the one-year rollout
 470 begins six months later. When the simulation is initialized on January 2, 2018, the heights
 471 in the arctic during boreal winter are approximately correct, while those in the antarctic
 472 are too cold (Figure 8d).

473 There is also a long-term drift toward lower heights in the subtropics and mid-latitudes,
 474 creating a roughly 30 m loss in Z_{500} by the end of the 1-year forecast.⁷ Climate models
 475 are tuned to avoid long-term drift in the predicted fields, but operational NWP models

⁷ 30 m deviation amounts to 0.5% of the full Z_{500} value and to 8.7% of the Z_{500} standard deviation (computed from the reanalysis data of the forecasted period).

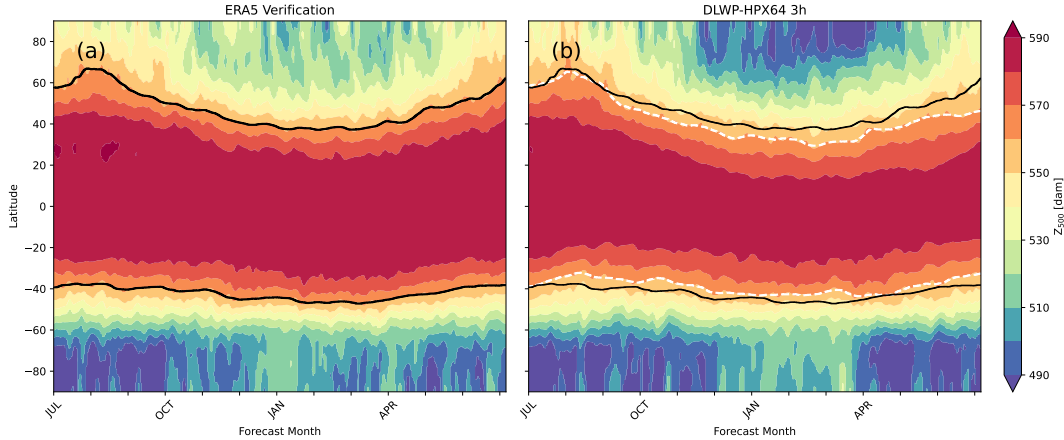


Figure 7: Zonally averaged three-day mean of Z_{500} plotted as a function of time and latitude for one year beginning on July 1 2017 for: (a) the ERA5 reanalysis, and (b) a recursive one-year rollout of the DLWP-HPX model. Also shown are 15-day averaged values of the 5600 m contour of Z_{500} for the ERA5 data (black lines) the DLWP-HPX simulation (white dashed lines).

476 are not so tuned. For example, significant model biases that grow over a time scale of
 477 several weeks are removed to create sub-seasonal ECMWF IFS S2S forecasts (Vitart, 2004;
 478 Weigel et al., 2008). To facilitate comparison of model drift with the ERA5 reanalysis,
 479 the pair of black lines in both panels show the 15-day mean of the zonally averaged 560-
 480 dam Z_{500} contours in the northern and southern hemisphere. The white lines in Figure 7b
 481 show the corresponding 560-dam Z_{500} contours for the DLWP-HPX simulation. The drift
 482 toward lower heights starts to become evident after two months in the northern hemi-
 483 sphere and continues to grow slowly for the remainder of the year. Differences show up
 484 earlier in the southern hemisphere, but the average drift is smaller and even disappears
 485 at a few times later in the year. As will be discussed in a forthcoming paper, both the
 486 errors near the poles and the drift in the tropics in Z_{500} can be corrected by incorporat-
 487 ing SST forecasts from a coupled atmosphere-ocean model.

488 The performance of three additional state-of-the-art DLWP models is compared
 489 with our model using this same metric in Figure 8, which shows the evolution of zonally
 490 averaged Z_{500} heights over a one-year rollout beginning January 2, 2018. This year is
 491 part of the test set for all of the models: our DLWP-HPX, Pangu-Weather, GraphCast,
 492 and FourCastNetv2 based on spherical Fourier neural operators (SFNO) (Bonev et al.,
 493 2023). Details about the code used to generate these rollouts can be found in section Ap-
 494 pendix B.

495 The Pangu-Weather model does not include solar forcing, and therefore, it does not
 496 follow the annual cycle. When stepped forward with a 24-h time step (Figure 8b), sig-
 497 nificant drift is apparent after about 1.5 months, which grows through the year without
 498 pushing the simulation into grossly unrealistic states. Based on the discussion of Extended
 499 Data, Fig. 7a in (Bi et al., 2023), one would not expect good performance from Pangu-
 500 Weather if rolled out with a 3-h time step, and indeed the 3-h rollout starts to produce
 501 significant errors after 1.5 months and generates completely unrealistic results after about
 502 5 months (Figure 8f). We nevertheless, show its performance to contrast it with our 3-
 503 h-time-resolution rollout (Figure 8e).

504 The version of GraphCast from NVIDIA’s Earth2MIP gives reasonable results for
 505 just the first 1.5 months (Figure 8c), while that from DeepMind goes bad after a cou-

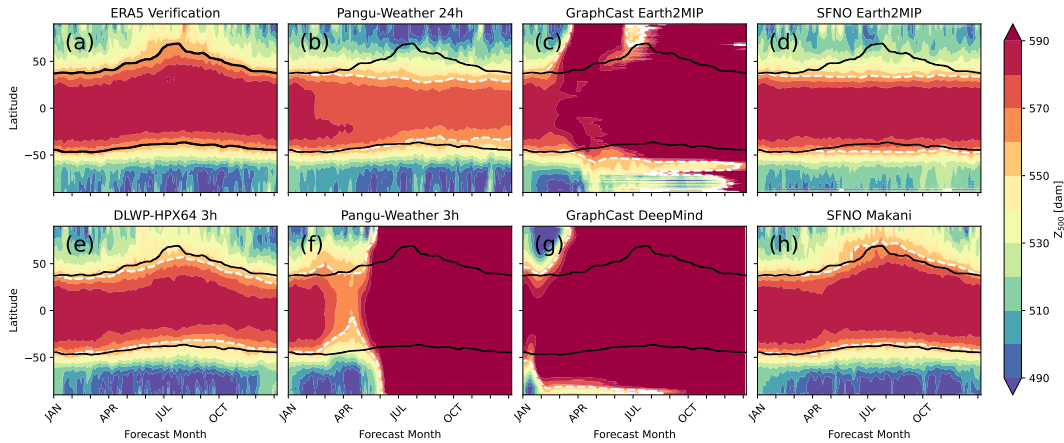


Figure 8: Zonally averaged three-day mean of Z_{500} plotted as a function of time and latitude: (a) for ERA5 reanalysis, (b)-(h) for recursive one-year simulations for each model as identified in the titles, initialized on January 2, 2018. Also shown are 15-day averaged values of the 5600 m contour of Z_{500} for the ERA5 data (black lines) each model simulation (white dashed lines).

506 ple weeks (Figure 8g). The SFNO Earth2MIP model (FourCastNetv2-small) shows es-
 507 sentially no drift over a full year (Figure 8d), although surprisingly, it does not follow
 508 the annual cycle despite including solar zenith angle as an input field. Some artifacts (hor-
 509 izontal stripes) are visible near the south pole within a month and at the north pole much
 510 later in the simulation. In contrast, the SFNO Makani model (Figure 8h), also includes
 511 solar zenith angle as an input field, and it does follow the annual cycle reasonably well.
 512 On balance, the performance of the SFNO Makani model is roughly similar to our DLWP-
 513 HPX model; it has larger errors near the poles, but less drift in the tropics.

514 In an ablation study (not shown), we investigated the effect of the top-of-atmosphere
 515 solar forcing input on the 365-day DLWP-HPX rollout by training a model that did not
 516 receive solar forcing input. In that case, the model still generated a stable forecast over
 517 the entire rollout period, but did not produce the full annual cycle. Interestingly, that
 518 simulation did roughly approximate the transition from summer into a perpetual autumn.

519 One qualitative way to appreciate the stable behavior of our one-year simulations
 520 is illustrated by comparing a 360.5 day simulation initialized on 1 April 2017 (with 6 h
 521 time steps and 3 h resolution) and the corresponding 27 March 2018 reanalysis in Fig-
 522 ure 9. The roughly one-year lead time is well beyond the limits of atmospheric predictabil-
 523 ity, so there is no reason to expect a close match between simulation and reanalysis. The
 524 360.5-day simulation time was chosen to display the simulated strong low-pressure center
 525 in the northeastern Pacific. The intensity of the system is typical for strong systems
 526 in our simulation, but about 40 m higher than the strongest systems periodically appear-
 527 ing in the ERA5 reanalysis. Lower-amplitude signals also appear in the Z_{1000} field, which
 528 is somewhat less than 50 m too low in the tropics. On balance, the overall character of
 529 this late-March weather pattern is quite plausible.

530 A more quantitative assessment of any tendency of our model to distort the atmo-
 531 spheric state by damping or amplifying mid-latitude perturbations at different wavelengths
 532 is provided by the plots of the Z_{500} power spectral density around 45°N in Figure 10.
 533 These spectra are averaged over 208 biweekly forecasts from the 2017-2018 test set; as
 534 such, the initial spectrum in black represents the average state of the atmosphere in the
 535 ERA5 reanalysis.

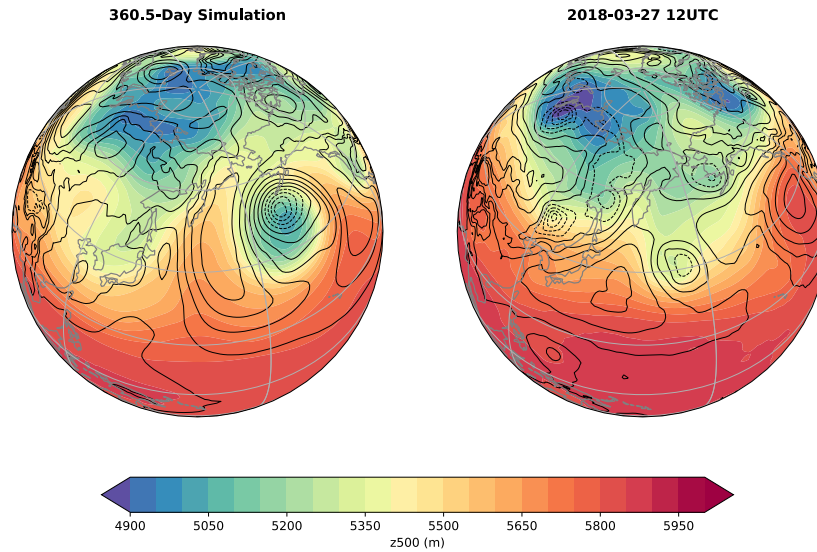


Figure 9: Z_{500} (color fill: 50 dam contour interval) and Z_{1000} (black contours: 40 m interval) for a free-running 360.5-day simulation and the corresponding ERA5 reanalysis for 00 UTC on 27 March 2018. Dashed black lines indicate values of $Z_{1000} \leq 40$ m (corresponding to sea-level pressures less than roughly 1008 hPa).

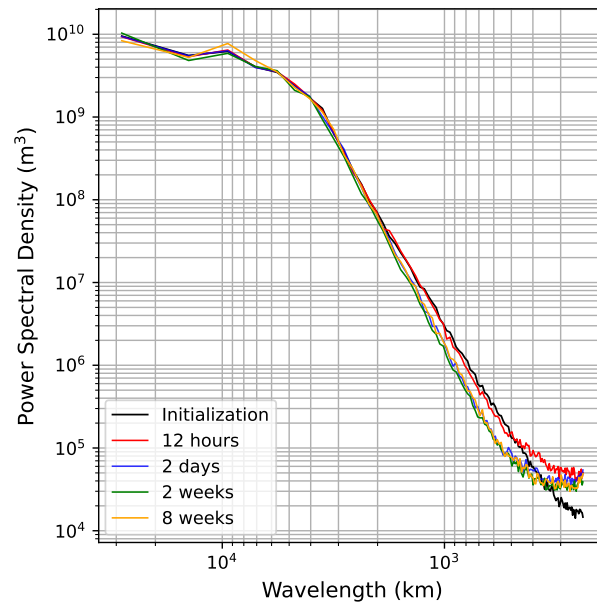


Figure 10: One dimensional power spectral density of the Z_{500} field around the 45° N latitude, averaged over 208 bi-weekly forecasts from 2017-2018 at: initialization (black), and at forecast lead times of 12 h, 2 d, 2, and 8 weeks.

Twelve hours (2 recursive steps) after initialization there is very little change in the spectra for wavelengths λ longer than 500 km (roughly 5 grid intervals), but the power in the shorter waves is amplified. Over the next 36 h, there is a gradual reduction in the amplitude at wavelengths $\lambda < 1800$ km to yield a spectrum that is modestly damped over the interval $380 < \lambda < 1800$ km and amplified at the shortest wavelengths. Surprisingly, the spectral distribution at two days remains essentially unchanged through at least sub-seasonal forecast lead times of eight weeks, which is consistent with the impression obtained examining images such as those in Figure 9. There is no gradual amplification or loss of amplitude in the simulated atmospheric systems after the first two days.

5 Conclusion

We have presented an improved CNN-based DLWP-HPX model that stably forecasts atmospheric evolution over a full one-year cycle using a very limited set of prognostic variables. The number of actual degrees of freedom characterising predictable atmospheric states at forecast lead times beyond 3–5 days is not known, but is far less than the total number of prognostic variables carried at every grid cell in state-of-the-art NWP models. Here, we have demonstrated that realistic atmospheric simulations can be performed using just seven prognostic variables above each node on a HEALPix mesh with 110 km between the nodes.

The HEALPix mesh (Gorski et al., 2005) has been used in astronomy for almost two decades, but has previously seen very little use in atmospheric science. The mesh covers the sphere with a hierarchical grid of equal-area cells uniformly spaced along circles at constant latitudes. A particularly important advantage of the HEALPix mesh for weather forecasting with CNNs is that it is an “east to the right” mesh, i.e., east has the same orientation in every HEALPix cell. Weather systems tend to travel west-to-east in mid- and high-latitudes and both east-to-west (tropical cyclones) or west-to-east (Madden-Julian Oscillation, convectively coupled Kelvin waves) in the tropics. The kernel weights in our convolutional stencils can more economically learn this behavior than on our previous cubed sphere mesh in which the eastward orientation across the stencil varies with longitude, particularly on the polar faces. Although switching from a cube-sphere mesh with 64×64 cells on each of the six faces to a HEALPix mesh with 32×32 cells on each of the 12 faces reduces the total number of grid points covering the sphere by half, it improves the Z_{500} RMSE error by almost one day at a 4-day forecast lead time (Figure 5).

Two other significant improvements to our model architecture were obtained by adding recursion via GRUs and by inverting the standard way channel depth is refined at deeper layers in the U-Net. In contrast to the original U-Net architecture Ronneberger et al. (2015), our channel depth halves instead of doubles as the spatial resolution is also halved in each successively deeper U-Net layer. This allows the model to devote more trainable parameters to describing the wide variety of fine-scale weather patterns while using comparatively fewer parameters to describe the simpler set of global weather patterns. Although this modification pushes the U-Net toward the basic ResNet architecture (He et al., 2016), we find the deeper U-Net layers continue to provide significant skill to the forecasts.

Additional modest improvements were implemented by switching to the GELU activation function and to 2×2 transposed strided convolutions when up-sampling; by increasing the total number of trainable parameters from 2.7 M to 9.8 M, adding the Z_{250} field, increasing the resolution to HPX64, and increasing the time resolution to 3 h (which gives us a 6 h time step). The benefits of 3-h time resolution were only realized when the model included the GRUs. The 3-h time resolution gives a good forecast of the daily cycle of surface temperature, and the model also learns the difference in the range of that cycle between regions of tropical forest and desert without geo-specific input data.

587 Finally, we replaced the pairs of successive convolutions in Weyn et al. (2020) with
 588 modified ConvNeXt blocks. The switch to the ConvNeXt blocks was only advantageous
 589 at higher resolutions, where in addition to improving accuracy, it reduced the memory
 590 footprint.

591 At one-week forecast lead time, the resulting model is roughly 1 day behind the
 592 ECMWF IFS S2S forecast error in Z_{500} RMSE and 1.5 days behind in ACC. These statis-
 593 tics are worse than those for Pangu-Weather (Bi et al., 2023) and GraphCast (Lam et
 594 al., 2022), both of which provide Z_{500} RMSE and ACC forecasts at $0.25^\circ \times 0.25^\circ$ reso-
 595 lution that are superior to the deterministic ECMWF IFS high-resolution model aver-
 596 aged to the same $0.25^\circ \times 0.25^\circ$ grid. Despite having less accuracy in medium range fore-
 597 casts, our model can be recursively stepped forward to generate better 500 hPa forecasts
 598 over seasonal and one-year rollouts than GraphCast and Pangu-Weather. It is also su-
 599 perior to the SFNO version of FourCastNetv2 currently on NVIDIA Earth2MIP, though
 600 it behaves similarly to the recently checkpointed version of SFNO Makani. Realistic low
 601 pressure systems and upper-level trough and ridge patterns continue to be generated by
 602 our model at the end of the one-year rollout.

603 Deep learning models for weather forecasting are evolving rapidly, with important
 604 advancements using a wide variety of architectures. Our DLWP-HPX model provides
 605 an example of what can be achieved using a relatively parsimonious approach. As such,
 606 it may be particularly useful for scientific investigations where it is advantageous to work
 607 with a minimal set of unknown variables to more concisely characterize sensitivities that
 608 might be revealed by techniques such as backpropagation with respect loss functions cus-
 609 tomized for analysis (as opposed to model training).

610 There are many avenues along which our DLPW-HPX model might be improved.
 611 One would be to adding additional prognostic fields while carefully examining the result-
 612 ing performance. Another one would lie in refining the CNN architecture, where the choice
 613 of particular inductive biases may be crucial (Thuemmel et al., 2023). A related impor-
 614 tant aspect of improving the modelled processes might be to incorporate explicit phys-
 615 ical constraints, yielding physics-informed differentiable artificial neural networks (Beucler
 616 et al., 2021; Shen et al., 2023). Other natural extensions of this work lie in examining
 617 the performance of the DLPW-HPX model in ensemble forecasts, which are crucial to
 618 sub-seasonal and seasonal prediction and to couple the atmospheric model with the ocean,
 619 thus moving toward a deep learning earth system model (Bauer et al., 2023). Prelimi-
 620 nary results suggest that coupling our model with a deep learning ocean model that pre-
 621 dictes sea surface temperatures (which are not incorporated in the current model) stabi-
 622 lizes the simulations and removes model drift in multi-decadal rollouts.

623 Appendix A Deep Learning on the HEALPix

624 A1 Seamless Evolution of Location Invariant Kernels

625 The Hierarchical Equal Area isoLatitude Pixelization (HEALPix) is a partitioning
 626 of the sphere that has found wide application in astronomy since it was introduced by
 627 Gorski et al. (2005). It divides the sphere into 12 base faces that can be hierarchically
 628 subdivided into patches of equal size. A key property for training CNNs on this mesh
 629 is the isolatitudinal alignment, that is, patches are aligned along lines of latitude and each
 630 patch has the same orientation, which we describe as “east to the right” in subsection 4.1.

631 To contrast and emphasize the difficulty that CNN kernels are facing on the cubed
 632 sphere mesh, we plot the lines of constant latitude on the six faces of the cubed sphere
 633 and on the twelve faces of the HEALPix in Figure A1. Except for the equator, all lines
 634 of constant latitude are bent on the cubed sphere, imposing challenges for a limited set
 635 of convolution kernels that must evolve location invariant pattern detectors and functions.
 636 For example, on the cubed sphere, kernels need to learn a wider range of behaviors to

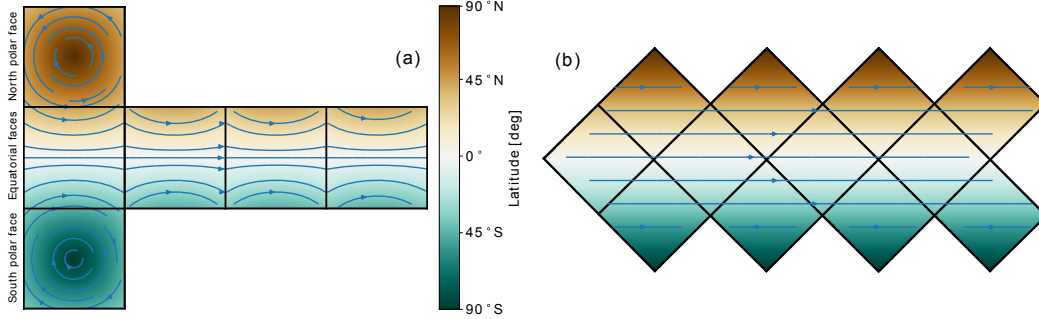


Figure A1: Lines of latitudes depicted as blue streamline arrows on the cubed sphere (a) and on the HEALPix (b). While the lines corresponding to constant eastward motion describe arcs of different radii on the cubed sphere mesh, the same motion translates to straight lines on the HEALPix mesh.

637 propagate eastward motions at the top-left versus the top right corners of the cubed sphere
 638 faces.

639 On the other hand, lines of constant latitude map to straight lines on the HEALPix
 640 mesh. This facilitates the formulation of location-invariant convolutional kernels for the
 641 propagation of weather systems, which tend to migrate eastward outside the tropics.

642 A2 Technical Implementation Details

643 Since deep learning libraries are optimized for image processing tasks, we consider
 644 each of the HEALPix’s 12 base faces as a regular two-dimensional tensor, i.e., we interpret
 645 the sphere as a composition of twelve images (cf. Figure 1 and Figure A2).

646 To simulate the spatial propagation of dynamics beyond individual faces, such that
 647 weather patterns can evolve globally on the sphere, we implement custom padding operations
 648 to concatenate the relevant information of all neighboring faces to each respective
 649 face of interest.

650 Figure A2 showcases our planet’s coastlines projected on the HEALPix faces in (a)
 651 and outlines the spatial organization of the twelve faces in (b). The arrangement of neighboring
 652 faces is exemplarily detailed for the northern (N) and southern (S) hemisphere,
 653 as well as for the equatorial faces (E). To simulate the neighborhood of, say, face E3, the
 654 face N2 must be concatenated to the left of E3, while face S3 is concatenated to the right.
 655 On the northern and southern hemispheres, neighboring faces are partially required to
 656 be rotated, as indicated in Figure A2 (c), (d), and (e).

657 A particular case occurs in the north and south corners of the tropical faces, where
 658 no natural neighbor exists—cf. Figure 1 and Figure A2 (f) for an illustration. To simulate
 659 the ninth neighbor of the respective corner, we interpolate the values from the according
 660 faces on the northern/southern hemisphere, by simply averaging the two corresponding
 661 values and writing the result in the simulated neighboring face. For example,
 662 to simulate the top left neighboring face of E3, we average the respective values from N2
 663 and N3, as detailed by the straight red arrows in Figure A2 (g). Values that do not lie
 664 on the main diagonal of the simulated face are not required to be interpolated, but are
 665 copied from the adjacent faces instead, denoted by the curved red arrows in Figure A2
 666 (g). The exemplary corner padding shows the case for the application of a 3×3 kernel
 667 with dilation of 1 or 2. Note that a 5×5 kernel could be applied in the same way. Importantly,
 668 the padding should not extend one neighboring face, which depends on the

669 resolution of the HEALPix mesh and the configuration of the applied convolution (ker-
 670 nel size and dilation). Otherwise, a hierarchy of padding operations would be required
 671 to be implemented and considered.

672 Appendix B DLWP Model Details

673 Configuration and parameter counts of all layers in our best performing model are
 674 detailed in Table B1, where c_{in} denotes the number of input channels, k is the kernel size,
 675 s the stride, and d the dilation. The receptive field of each layer with respect to the net-
 676 work input is reported under RF and the output shape takes (F, H, W, C) with F the
 677 number of faces, H and W height and width, and C the number of output channels. The
 678 dashed line separates the model’s encoder (above) and decoder (below) components. All
 679 ConvNeXt- and GRU-blocks are additionally broken down into their operations, visualized
 680 by the indented layer names. Numbers in brackets following individual layer names cor-
 681 respond to outputs, which are concatenated to the respective Concat layers in the de-
 682 coder. All convolution layers with $k = 3$ are followed by GELU activation functions.
 683 Residual connections are not reported as they neither change the spatial resolution nor
 684 the number of channels, and they do not contribute to the parameter count. Color codes
 685 approximate those used in the model schematic in Figure 2.

686 To generate 1-year rollouts for Pangu-Weather, GraphCast, and FourCastNet2 (SFNO),
 687 as plotted in Figure 8, we considered the respective public repositories with the pretrained
 688 model weights. More concretely, we generated the SFNO Earth2MIP (`fcnv2_sm`) and
 689 GraphCast Earth2MIP (`graphcast`) forecasts with NVIDIA’s `earth2mip` package,⁸ specif-
 690 ically developing a custom script for long rollouts.⁹ The SFNO Makani forecast, which
 691 responds reasonably to the solar forcing by receiving an additional $\cos \phi$ input (where
 692 ϕ is the solar zenith angle), was generated with NVIDIA’s Makani package.¹⁰ Interest-
 693 ingly, the original GraphCast DeepMind code base¹¹ produced slightly different results
 694 and saturated even faster than the Earth2MIP version, which might result from differ-
 695 ent random seeds. For the DeepMind version of GraphCast, we downloaded the model
 696 weights¹² provided through their repository. Pangu-Weather forecasts in 24 h and 3 h res-
 697 olution (with respective checkpoint files for the 24 h¹³ and 3 h¹⁴ models) were generated
 698 by using the original repository.¹⁵

699 Open Research Section

700 Instructions for training, and a trained model for inference, are available at [https://](https://github.com/CognitiveModeling/dlwp-hpx/)
 701 github.com/CognitiveModeling/dlwp-hpx/. In addition, PyTorch code for training
 702 the DLWP-HPX models along with checkpoints of trained models will be provided via
 703 NVIDIA’s Modulus framework. Accompanying scripts for data preprocessing, including
 704 the projection to and from the HEALPix mesh, as well as postprocessing utilities, includ-
 705 ing evaluation routines, will be made available in the repository at [https://github.com/](https://github.com/NVIDIA/modulus/tree/main/examples/weather)
 706 [NVIDIA/modulus/tree/main/examples/weather](https://github.com/NVIDIA/modulus/tree/main/examples/weather). All spherical shells of data from ERA5
 707 (Hersbach et al., 2020) were downloaded from Copernicus, where variables on various con-

⁸ <https://github.com/NVIDIA/earth2mip>

⁹ https://github.com/NVIDIA/earth2mip/blob/main/examples/utils/workflows/1_year_run.py

¹⁰ <https://github.com/NVIDIA/makani>

¹¹ <https://github.com/google-deepmind/graphcast>

¹² https://storage.googleapis.com/dm_graphcast/params/GraphCast%20-%20ERA5%201979-2017%20-%20resolution%200.25%20-%20pressure%20levels%2037%20-%20mesh%20to6%20-%20precipitation%20input%20and%20output.npz

¹³ <https://drive.google.com/file/d/1lweQ1xcn9fG0zKNW8ne1Khr9ehRTI6HP/view>

¹⁴ <https://drive.google.com/file/d/1EdoLlAXqE9iZLt9Ej9i-JW9LTJ9Jtewt/view>

¹⁵ <https://github.com/198808xc/Pangu-Weather>

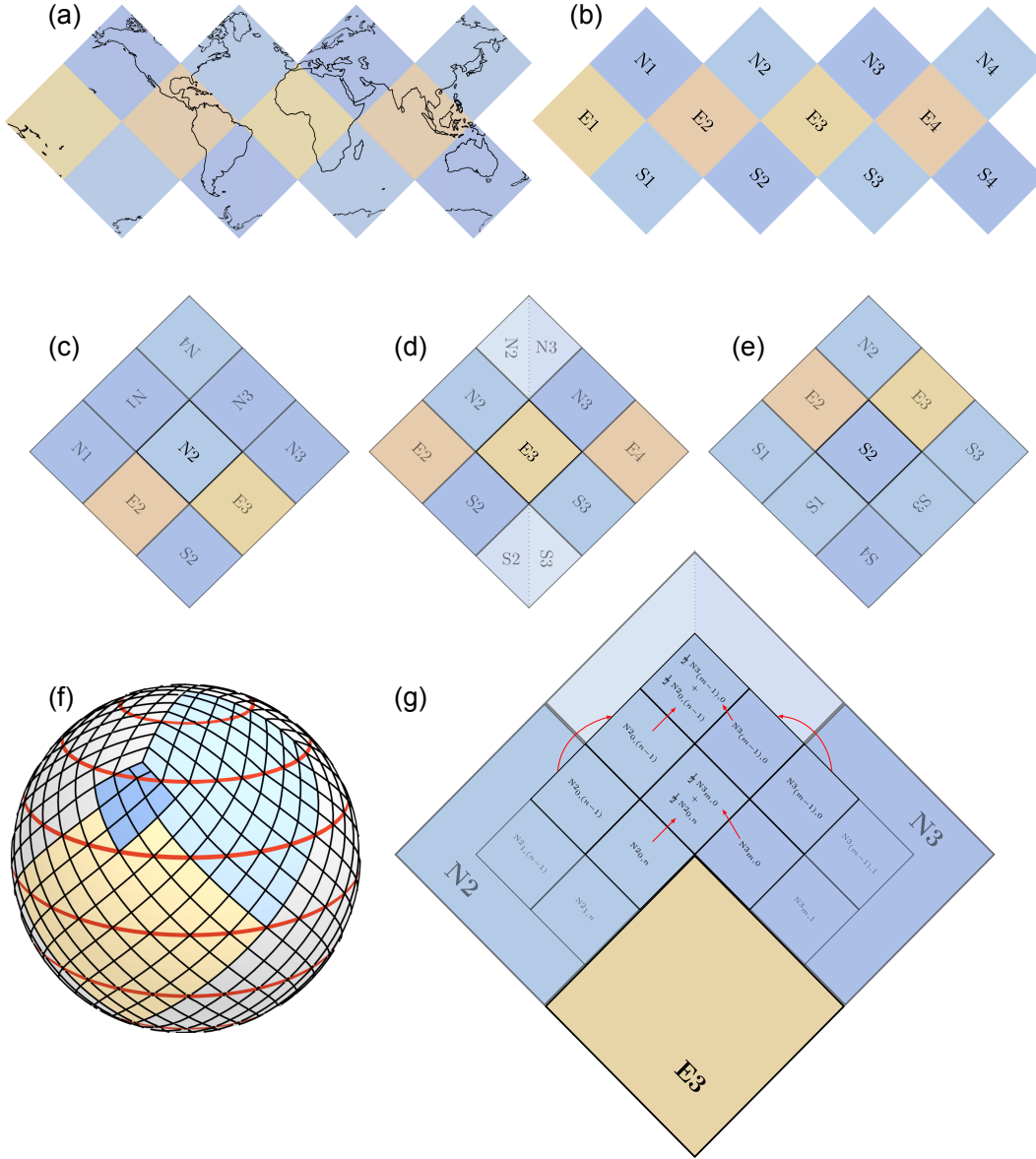


Figure A2: 2D HEALPix face arrangement and padding. (a) depicts the distribution of coastlines over the twelve HEALPix faces. (b) enumerates the twelve faces of the HEALPix with each four faces on the northern and southern hemisphere and around the equator. (c), (d), and (e): Exemplary alignment and rotations of neighboring faces before applying the padding operation on northern (c), equatorial (d), and southern faces (e). (f) emphasizes the special corner case, which is detailed in (g) to visualize the padding, where a ninth pixel is simulated by averaging the two respective values from the adjacent faces.

Table B1: Details of the best performing model. Description of color codes and abbreviations are reported in section Appendix B

Layer	c_{in}	k	s	d	RF	Output shape	Parameter count		
							Weights	Biases	Σ
ConvNeXt									
Conv2d	18	1	1	1	1×1	(12, 64, 64, 136)	2 448	136	2 584
Conv2d	18	3	1	1	3×3	(12, 64, 64, 544)	88 128	544	88 672
Conv2d	544	3	1	1	5×5	(12, 64, 64, 544)	2 663 424	544	2 663 968
Conv2d (1)	544	1	1	1	5×5	(12, 64, 64, 136)	73 984	136	74 120
AvgPool2d	136	2	2	—	6×6	(12, 32, 32, 136)	0	0	0
ConvNeXt									
Conv2d	136	1	1	1	6×6	(12, 32, 32, 68)	9 248	68	9 316
Conv2d	136	3	1	2	14×14	(12, 32, 32, 272)	332 928	272	333 200
Conv2d	272	3	1	2	22×22	(12, 32, 32, 272)	665 856	272	666 128
Conv2d (2)	272	1	1	1	22×22	(12, 32, 32, 68)	18 496	68	18 564
AvgPool2d	68	2	2	—	24×24	(12, 16, 16, 68)	0	0	0
ConvNeXt									
Conv2d	68	1	1	1	24×24	(12, 16, 16, 34)	2 312	34	2 346
Conv2d	68	3	1	4	56×56	(12, 16, 16, 136)	83 232	136	83 368
Conv2d	136	3	1	4	88×88	(12, 16, 16, 136)	166 464	136	166 600
Conv2d	136	1	1	1	88×88	(12, 16, 16, 34)	4 624	34	4 658
ConvNeXt									
Conv2d	34	1	1	1	88×88	(12, 16, 16, 68)	2 312	68	2 380
Conv2d	34	3	1	4	120×120	(12, 16, 16, 136)	41 616	136	41 752
Conv2d	136	3	1	4	152×152	(12, 16, 16, 136)	166 464	136	166 600
Conv2d	136	1	1	1	152×152	(12, 16, 16, 68)	9 248	68	9 316
GRU									
Conv2d	136	1	1	1	152×152	(12, 16, 16, 136)	18 496	136	18 632
Conv2d	136	1	1	1	152×152	(12, 16, 16, 68)	9 248	68	9 316
ConvTrans2d	68	2	2	1	154×154	(12, 32, 32, 68)	18 496	68	18 476
Concat (2)	—	—	—	—	—	(12, 32, 32, 136)	0	0	0
ConvNeXt									
Conv2d	136	3	1	2	154×154	(12, 32, 32, 272)	332 928	272	333 200
Conv2d	272	3	1	2	162×162	(12, 32, 32, 272)	665 856	272	666 128
Conv2d	272	1	1	1	170×170	(12, 32, 32, 136)	36 992	136	37 128
GRU									
Conv2d	272	1	1	1	170×170	(12, 32, 32, 272)	73 984	272	74 256
Conv2d	136	1	1	1	170×170	(12, 32, 32, 136)	36 992	136	37 128
ConvTrans2d	136	2	2	1	171×171	(12, 64, 64, 136)	73 984	136	74 120
Concat (1)	—	—	—	—	—	(12, 64, 64, 272)	0	0	0
ConvNeXt									
Conv2d	272	1	1	1	171×171	(12, 64, 64, 136)	36 992	136	37 128
Conv2d	272	3	1	1	173×173	(12, 64, 64, 544)	1 331 712	544	1 332 256
Conv2d	544	3	1	1	175×175	(12, 64, 64, 544)	2 663 424	544	2 663 968
Conv2d	544	1	1	1	175×175	(12, 64, 64, 136)	73 984	136	74 120
GRU									
Conv2d	272	1	1	1	175×175	(12, 64, 64, 272)	73 984	272	74 256
Conv2d	272	1	1	1	175×175	(12, 64, 64, 136)	36 992	136	37 128
Conv2d	136	1	1	1	175×175	(12, 64, 64, 14)	1 904	14	1 918
							9 816 752	6 066	9 822 818

stant pressure levels, such as Z_{500} or T_{850} , and variables on single levels, such as T_{2m} or $TCWV$, are hosted open to the public, available at <https://cds.climate.copernicus.eu/cdsapp#!/dataset/reanalysis-era5-pressure-levels?tab=form> and <https://cds.climate.copernicus.eu/cdsapp#!/dataset/reanalysis-era5-single-levels?tab=overview>.

Acknowledgments

We would like to thank Mauro Bisson from NVIDIA Corp. for providing optimized CUDA kernels for the HEALPix padding implementation, and Jonathan Weyn who previously implemented a code base on which this work was built. We thank Peter Düben and a second anonymous reviewer for encouraging us to generate and compare the 1-year rollouts for other state-of-the-art DLWP methods and for other valuable suggestions. This work received funding from Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) under Germany’s Excellence Strategy EXC 2064 – 390727645 and from the Office of Naval Research under grants N0014-21-1-2827 and N00014-22-1-2807. We thank the Deutscher Akademischer Austauschdienst (DAAD, German Academic Exchange Service) as well as the International Max Planck Research School for Intelligent Systems (IMPRS-IS) for supporting Matthias Karlbauer. Nathaniel was supported by a National Defense Science and Engineering Graduate Fellowship. We are grateful to NVIDIA and Stan Posey for the donation of A100 GPU cards. This research was additionally supported by a grant from the NVIDIA Applied Research Accelerator Program and utilized an NVIDIA DGX-100 Workstation. Moreover, this work benefited substantially from the barrier-free high quality ERA5 dataset provided by the ECMWF.

Author Roles

Matthias implemented model, training and evaluation routines in PyTorch, as well as the HEALPix-related projection scripts under consideration of the `healpy` package, and drafted the manuscript together with Dale who supervised this project closely and who also made the model schematic in Figure 2. Nathaniel was involved in discussions about model evolution and code structures and generated Figure 6, Figure 7, and Figure 10. Raul was involved in model discussions and generated Figure 9. Thorsten helped with implementing the distributed PyTorch pipeline for multi-GPU training and with accelerating the process pipeline. Noah Brenowitz and Boris Bonev generated the 365-days rollouts with the Earth2MIP and Makani packages for SFNO and GraphCast. Martin co-supervised this project and helped with proofreading and writing.

References

- Battaglia, P. W., Hamrick, J. B., Bapst, V., Sanchez-Gonzalez, A., Zambaldi, V., Malinowski, M., . . . others (2018). Relational inductive biases, deep learning, and graph networks. *arXiv preprint arXiv:1806.01261*.
- Bauer, P., Dueben, P., Chantry, M., Doblus-Reyes, F., Hoefler, T., McGovern, A., & Stevens, B. (2023). Deep learning and a changing economy in weather and climate prediction. *Nature Reviews Earth & Environment*, 4(8), 507–509. Retrieved from <https://doi.org/10.1038/s43017-023-00468-z> doi: 10.1038/s43017-023-00468-z
- Bauer, P., Thorpe, A., & Brunet, G. (2015). The quiet revolution of numerical weather prediction. *Nature*, 525(7567), 47–55.
- Benjamin, S. G., Brown, J. M., Brunet, G., Lynch, P., Saito, K., & Schlatter, T. W. (2019). 100 years of progress in forecasting and nwp applications. *Meteorological Monographs*, 59, 13–1.
- Beucler, T., Pritchard, M., Rasp, S., Ott, J., Baldi, P., & Gentine, P. (2021). Enforcing analytic constraints in neural networks emulating physical systems. *Physical*

- 757 *Review Letters*, 126(9), 098302.
- 758 Bi, K., Xie, L., Zhang, H., Chen, X., Gu, X., & Tian, Q. (2023). Accurate medium-
759 range global weather forecasting with 3d neural networks. *Nature*. doi: doi.org/
760 10.1038/s41586-023-06185-3
- 761 Bonev, B., Kurth, T., Hundt, C., Pathak, J., Baust, M., Kashinath, K., & Anandku-
762 mar, A. (2023). Spherical fourier neural operators: Learning stable dynamics on
763 the sphere. *arXiv preprint arXiv:2306.03838*.
- 764 Charney, J. G., Fjörtoft, R., & Neumann, J. V. (1950). Numerical Integration of the
765 Barotropic Vorticity Equation. *Tellus A*, 2(4).
- 766 Chen, K., Han, T., Gong, J., Bai, L., Ling, F., Luo, J.-J., . . . Ouyang, W. (2023).
767 Fengwu: Pushing the skillful global medium-range weather forecast beyond 10
768 days lead. *arXiv preprint arXiv:2304.02948*.
- 769 Cho, K., van Merriënboer, B., Gulcehre, C., Bougares, F., Schwenk, H., & Bengio, Y.
770 (2014). Learning phrase representations using rnn encoder-decoder for statistical
771 machine translation. In *Conference on empirical methods in natural language
772 processing (emnlp 2014)*.
- 773 Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner,
774 T., . . . others (2020). An image is worth 16x16 words: Transformers for image
775 recognition at scale. *arXiv preprint arXiv:2010.11929*.
- 776 Dueben, P. D., & Bauer, P. (2018). Challenges and design choices for global weather
777 and climate models based on machine learning. *Geoscientific Model Develop-*
778 *ment*, 11(10), 3999–4009.
- 779 Gori, M., Monfardini, G., & Scarselli, F. (2005). A new model for learning in graph
780 domains. In *Proceedings. 2005 IEEE International Joint Conference on Neural Net-*
781 *works, 2005*. (Vol. 2, pp. 729–734).
- 782 Gorski, K. M., Hivon, E., Banday, A. J., Wandelt, B. D., Hansen, F. K., Reinecke,
783 M., & Bartelmann, M. (2005). Healpix: A framework for high-resolution
784 discretization and fast analysis of data distributed on the sphere. *The Astro-*
785 *physical Journal*, 622(2), 759.
- 786 Guibas, J., Mardani, M., Li, Z., Tao, A., Anandkumar, A., & Catanzaro, B. (2021).
787 Efficient token mixing for transformers via adaptive fourier neural operators. In
788 *International conference on learning representations*.
- 789 He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep residual learning for image recog-
790 nition. In *Proceedings of the IEEE conference on computer vision and pattern
791 recognition* (pp. 770–778).
- 792 Hendrycks, D., & Gimpel, K. (2016). Gaussian error linear units (gelus). *arXiv
793 preprint arXiv:1606.08415*.
- 794 Hersbach, H., Bell, B., Berrisford, P., Hirahara, S., Horányi, A., Muñoz-Sabater, J.,
795 . . . others (2020). The era5 global reanalysis. *Quarterly Journal of the Royal
796 Meteorological Society*, 146(730), 1999–2049.
- 797 Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. *Neural computa-*
798 *tion*, 9(8), 1735–1780.
- 799 Hu, Y., Chen, L., Wang, Z., & Li, H. (2022). Swinvrnn: A data-driven ensemble
800 forecasting model via learned distribution perturbation. *arXiv preprint
801 arXiv:2205.13158*.
- 802 Huang, H., Lin, L., Tong, R., Hu, H., Zhang, Q., Iwamoto, Y., . . . Wu, J. (2020).
803 Unet 3+: A full-scale connected unet for medical image segmentation. In
804 *Icassp 2020-2020 IEEE International Conference on Acoustics, Speech and Signal
805 Processing (ICASSP)* (pp. 1055–1059).
- 806 Keisler, R. (2022). Forecasting global weather with graph neural networks. *arXiv
807 preprint arXiv:2202.07575*.
- 808 Kingma, D. P., & Ba, J. (2014). Adam: A method for stochastic optimization. *arXiv
809 preprint arXiv:1412.6980*.
- 810 Kipf, T. N., & Welling, M. (2016). Semi-supervised classification with graph convolu-
811 tional networks. *arXiv preprint arXiv:1609.02907*.

- 812 Krachmalnicoff, N., & Tomasi, M. (2019). Convolutional neural networks on
 813 the healpix sphere: a pixel-based algorithm and its application to cmb data
 814 analysis. *Astronomy & Astrophysics*, *628*, A129.
- 815 Kurth, T., Subramanian, S., Harrington, P., Pathak, J., Mardani, M., Hall, D., ...
 816 Anandkumar, A. (2022). Fourcastnet: Accelerating global high-resolution
 817 weather forecasting using adaptive fourier neural operators. *arXiv preprint*
 818 *arXiv:2208.05419*.
- 819 Lam, R., Sanchez-Gonzalez, A., Willson, M., Wirnsberger, P., Fortunato, M., Pritzel,
 820 A., ... others (2022). Graphcast: Learning skillful medium-range global weather
 821 forecasting. *arXiv preprint arXiv:2212.12794*.
- 822 Li, Z., Kovachki, N., Aizzadenesheli, K., Liu, B., Bhattacharya, K., Stuart, A.,
 823 & Anandkumar, A. (2020). Fourier neural operator for parametric partial
 824 differential equations. *arXiv preprint arXiv:2010.08895*.
- 825 Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., ... Guo, B. (2021). Swin trans-
 826 former: Hierarchical vision transformer using shifted windows. In *Proceedings of*
 827 *the ieee/cvf international conference on computer vision* (pp. 10012–10022).
- 828 Liu, Z., Mao, H., Wu, C.-Y., Feichtenhofer, C., Darrell, T., & Xie, S. (2022). A con-
 829 vnet for the 2020s. In *Proceedings of the ieee/cvf conference on computer vision*
 830 *and pattern recognition* (pp. 11976–11986).
- 831 Lopez-Gomez, I., McGovern, A., Agrawal, S., & Hickey, J. (2022). Global extreme
 832 heat forecasting using neural weather models. *arXiv preprint arXiv:2205.10972*.
- 833 Lorenz, E. N. (1969). The predictability of a flow which possesses many scales of mo-
 834 tion. *Tellus*, *21*(3), 289–307.
- 835 Loshchilov, I., & Hutter, F. (2016). Sgdr: Stochastic gradient descent with warm
 836 restarts. In *International conference on learning representations*.
- 837 Palmer, T. (2019). The ecmwf ensemble prediction system: Looking back (more than)
 838 25 years and projecting forward 25 years. *Quarterly Journal of the Royal Mete-
 839 orological Society*, *145*, 12–24.
- 840 Pathak, J., Subramanian, S., Harrington, P., Raja, S., Chattopadhyay, A., Mardani,
 841 M., ... others (2022). Fourcastnet: A global data-driven high-resolution
 842 weather model using adaptive fourier neural operators. *arXiv preprint*
 843 *arXiv:2202.11214*.
- 844 Perraudin, N., Defferrard, M., Kacprzak, T., & Sgier, R. (2019). DeepSphere: Efficient
 845 spherical convolutional neural network with healpix sampling for cosmological
 846 applications. *Astronomy and Computing*, *27*, 130–146.
- 847 Pfaff, T., Fortunato, M., Sanchez-Gonzalez, A., & Battaglia, P. W. (2020). Learning
 848 mesh-based simulation with graph networks. *arXiv preprint arXiv:2010.03409*.
- 849 Rasp, S., & Thuerey, N. (2021). Data-driven medium-range weather prediction with a
 850 resnet pretrained on climate simulations: A new model for weatherbench. *Jour-
 851 nal of Advances in Modeling Earth Systems*, *13*(2), e2020MS002405.
- 852 Ronneberger, O., Fischer, P., & Brox, T. (2015). U-net: Convolutional networks for
 853 biomedical image segmentation. In *International conference on medical image*
 854 *computing and computer-assisted intervention* (pp. 234–241).
- 855 Scarselli, F., Gori, M., Tsoi, A. C., Hagenbuchner, M., & Monfardini, G. (2008). The
 856 graph neural network model. *IEEE transactions on neural networks*, *20*(1), 61–
 857 80.
- 858 Scher, S., & Messori, G. (2018). Predicting weather forecast uncertainty with machine
 859 learning. *Quarterly Journal of the Royal Meteorological Society*, *144*(717), 2830–
 860 2841.
- 861 Scher, S., & Messori, G. (2019). Weather and climate forecasting with neural net-
 862 works: using GCMs with different complexity as study-ground. *Geoscientific*
 863 *Model Development*, *12*, 2797–2809.
- 864 Shen, C., Appling, A. P., Gentine, P., Bandai, T., Gupta, H., Tartakovsky, A., ...
 865 Lawson, K. (2023). Differentiable modelling to unify machine learning and
 866 physical models for geosciences. *Nature Reviews Earth & Environment*, *4*(8),

- 552–567. Retrieved from <https://doi.org/10.1038/s43017-023-00450-9>
 doi: 10.1038/s43017-023-00450-9
- Thuemmel, J., Karlbauer, M., Otte, S., Zarfl, C., Martius, G., Ludwig, N., . . . others
 (2023). Inductive biases in deep learning models for weather prediction. *arXiv
 preprint arXiv:2304.04664*.
- Tobler, W. R. (1970). A computer movie simulating urban growth in the detroit re-
 gion. *Economic geography*, 46(sup1), 234–240.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., . . .
 Polosukhin, I. (2017). Attention is all you need. *Advances in neural information
 processing systems*, 30.
- Vitart, F. (2004). Monthly forecasting at ECMWF. *Monthly Weather Review*, 132,
 2761–2779. doi: 10.1175/MWR2826.1
- Weigel, A. P., Baggenstos, D., Liniger, M. A., Vitart, F., & Appenzeller, C. (2008).
 Probabilistic Verification of Monthly Temperature Forecasts. *Monthly Weather
 Review*, 136, 5162–5182. doi: 10.1175/2008MWR2551.1
- Weyn, J. A., Durran, D. R., & Caruana, R. (2019). Can machines learn to predict
 weather? using deep learning to predict gridded 500-hpa geopotential height
 from historical weather data. *Journal of Advances in Modeling Earth Systems*,
 11(8), 2680–2693.
- Weyn, J. A., Durran, D. R., & Caruana, R. (2020). Improving data-driven global
 weather prediction using deep convolutional neural networks on a cubed sphere.
Journal of Advances in Modeling Earth Systems, 12(9), e2020MS002109.
- Weyn, J. A., Durran, D. R., Caruana, R., & Cresswell-Clay, N. (2021). Sub-seasonal
 forecasting with a large ensemble of deep-learning weather prediction models.
Journal of Advances in Modeling Earth Systems, 13(7), e2021MS002502.
- Zhou, Z., Rahman Siddiquee, M. M., Tajbakhsh, N., & Liang, J. (2018). Unet++:
 A nested u-net architecture for medical image segmentation. In *Deep learning
 in medical image analysis and multimodal learning for clinical decision support*
 (pp. 3–11). Springer.