# Advancing Parsimonious Deep Learning Weather Prediction using the HEALPix Mesh

Matthias Karlbauer<sup>1,2</sup>, Nathaniel Cresswell-Clay<sup>2</sup>, Dale R Durran<sup>2</sup>, Raul A Moreno<sup>2</sup>, Thorsten Kurth<sup>3</sup>, and Martin V Butz<sup>1</sup>

<sup>1</sup>Neuro-Cognitive Modeling Group, University of Tübingen <sup>2</sup>Department of Atmospheric Sciences, University of Washington <sup>3</sup>NVIDIA Switzerland AG

March 13, 2024

# Advancing Parsimonious Deep Learning Weather Prediction using the HEALPix Mesh

### Matthias Karlbauer<sup>1</sup>, Nathaniel Cresswell-Clay<sup>2</sup>, Dale R. Durran<sup>2</sup>, Raul A. Moreno<sup>2</sup>, Thorsten Kurth<sup>3</sup>, Boris Bonev<sup>3</sup>, Noah Brenowitz<sup>4</sup>, and Martin V. Butz<sup>1</sup>

6	$^{1}\mathrm{Neuro-Cognitive}$ Modeling Group, Department of Computer Science, University of Tübingen, Tübingen,
7	Germany
8	<sup>2</sup> Department of Atmospheric Sciences, University of Washington, Seattle, WA, USA
9	<sup>3</sup> NVIDIA Switzerland AG, Zürich, Switzerland
10	$^{4}$ NVIDIA Corporation, Seattle, USA

Key Points:

1

2

3

4

5

11

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

Corresponding author: Dale R. Durran, drdee@uw.edu

#### 18 Abstract

We present a parsimonious deep learning weather prediction model on the Hierarchical 19 Equal Area isoLatitude Pixelization (HEALPix) to forecast seven atmospheric variables 20 for arbitrarily long lead times on a global approximately 110 km mesh at 3h time res-21 olution. In comparison to state-of-the-art machine learning weather forecast models, such 22 as Pangu-Weather and GraphCast, our DLWP-HPX model uses coarser resolution and 23 far fewer prognostic variables. Yet, at one-week lead times its skill is only about one day 24 behind the state-of-the-art numerical weather prediction model from the European Cen-25 tre for Medium-Range Weather Forecasts. We report successive forecast improvements 26 resulting from model design and data-related decisions, such as switching from the cubed 27 sphere to the HEALPix mesh, inverting the channel depth of the U-Net, and introduc-28 ing gated recurrent units (GRU) on each level of the U-Net hierarchy. The consistent 20 east-west orientation of all cells on the HEALPix mesh facilitates the development of location-30 invariant convolution kernels that are successfully applied to propagate global weather 31 patterns across our planet. Without any loss of spectral power after two days, the model 32 can be unrolled autoregressively for hundreds of steps into the future to generate stable 33 and realistic states of the atmosphere that respect seasonal trends, as showcased in one-34 year simulations. Our parsimonious DLWP-HPX model is research-friendly and poten-35

tially well-suited for sub-seasonal and seasonal forecasting.

#### <sup>37</sup> Plain Language Summary

Weather forecasting is traditionally realized by numerical weather prediction mod-38 els that solve physical equations to simulate the progression of the atmosphere. Numer-39 ical methods are compute intense and their performance is increasingly challenged by 40 less compute demanding but still highly sophisticated machine learning approaches. Yet, 41 a downside of these new models is their reliability: They are not guaranteed to gener-42 ate physically plausible states, which often prevents them from generating stable and re-43 alistic forecasts beyond two weeks into the future. Here, a parsimonious machine learn-44 ing 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 mod-46 els) over an entire year. Despite the small number of variables, our model generates fore-47 casts that only fall behind expensive state-of-the-art predictions by a single day. That 48 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 learn-50 ing models beyond two weeks promises to extend horizons for planning in various fields 51 that impact environment, economy, and society. 52

#### 53 1 Introduction

Four years ago, Weyn et al. (2019) posed the question "Can machines learn to pre-54 dict the weather?" and demonstrated that data driven convolutional neural networks can 55 forecast the evolution of the 500 hPa surface much better than the alternative dynam-56 ical model, the barotropic vorticity equation, which was used in the first numerical weather 57 prediction (NWP) model (Charney et al., 1950). An extremely rapid evolution of deep 58 learning weather prediction (DLWP) models followed, culminating in the recent Pangu-59 Weather (Bi et al., 2023) and GraphCast models (Lam et al., 2022), which outperform 60 the deterministic forecast from the state-of-the-art Integrated Forecast System (IFS) of 61 the European Centre for Medium-Range Weather Forecasts (ECMWF). 62

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

forecasts is enormous and only available at a handful of advanced dedicated centers. Ensemble forecasts, which provide an important way to account for uncertainty by generating a set of equally plausible predictions and extend the limit of skillful forecasts beyond that of a single deterministic model run, are also limited by the computational burden of high-resolution NWP to about 50 members (Palmer, 2019).

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

The actual number of independent degrees of freedom required to represent the pre-86 dictable components of the global atmosphere is unknown, but it clearly decreases with 87 increasing forecast lead times (Lorenz, 1969). GraphCast (Lam et al., 2022), for exam-88 ple, has achieved success at lead times as short as 6 h with 227 spherical shells of data. It can produce forecasts using much less computation time than the ECMWF IFS, but 90 it still requires large computing resources for training: 3 weeks using 32 TPU 4 proces-91 sors. Pangu-Weather (Bi et al., 2023) cuts the number of spherical shells by almost 2/392 to 69. The spherical Fourier neural operator (SFNO) version of FourCastNet compared 93 with the IFS in (Bonev et al., 2023) uses 73 spherical shells of data. Here, we take this 94 reduction much farther, presenting a parsimonious DLWP model that uses just 7 spher-95 ical shells of data to efficiently provide forecasts approaching the skill of ECMWF. While 96 not as accurate as GraphCast or Pangu-Weather for medium range forecasts with lead 97 times less than two weeks, we demonstrate that our model generates far less bias in fore-98 casts of 500 hPa height in one-year iterative forecasts. In addition, our model is poten-99 tially better suited for research applications such as computing the sensitivities of its com-100 pact state vector to custom diagnostic functions by backpropagation. 101

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

#### <sup>110</sup> 2 Related Work

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

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

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

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

In contrast to the preceding approaches, graph neural networks (Gori et al., 2005;
Scarselli et al., 2008; Kipf & Welling, 2016; Battaglia et al., 2018; Pfaff et al., 2020) where
applied on icosahedral meshes at course resolution by Keisler (2022) and at fine resolution in the Graphcast model (Lam et al., 2022). Similarly to Pangu-Weather, GraphCast
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

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

<sup>&</sup>lt;sup>1</sup> The 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

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

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

The grid spacing, or shortest inter-node spacing, on the HEALPix mesh is the diagonal distance between a pair of nodes on adjacent latitude lines. Denoting a HEALPix mesh with *n* divisions along one side of the original 12 faces as HPX*n*. The grid spacing is approximately 220 km ( $\approx 2^{\circ}$ ) for HPX32 and 110 km ( $\approx 1^{\circ}$ ) for HPX64.<sup>3</sup>

#### <sup>193</sup> 3.2 Machine Learning Architecture

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

202

# 3.2.1 Residual Prediction

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

<sup>&</sup>lt;sup>2</sup> https://github.com/CognitiveModeling/dlwp-hpx

 $<sup>^{3}</sup>$ 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.

 $<sup>^{4}</sup>$  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  $3 \times 3$  stencil, as applied by a standard convolution. To apply the  $3 \times 3$  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.

differential equations, and has been used successfully in previous deep-learning weather 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

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

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

#### 227 3.2.3 Recurrent Modules

The vanilla U-Net is a feed-forward network, which treats successive inputs independently even if the data represents a continuous sequence over time. Feed-forward networks do not have any memory capacity. They do not maintain an internal state between time steps. To enable the exploitation of information from previous latent states, we include 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 con guration 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.

with kernel size k = 1. We chose GRUs over LSTMs(Hochreiter & Schmidhuber, 1997)
 since we re-initialize the recurrent data over each24 h-cycle, and therefore do not require
 forget-gates (as con rmed experimentally, not shown).

3.2.4 Miscellaneous Modi cations

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

#### 3.2.5 Time Stepping Scheme

244

Similarly to Weyn et al. (2021), we apply a two-in-two-out mapping with a temporal resolution twice as ne as the actual time step. For example, two atmospheric states 3 h apart (each consisting of seven prognostic, along with three prescribed elds) are concatenated and input to the model, which generates a new pair of states, each characterising the atmosphere6 h later in time. This strategy is observed to stabilize and accelerate the training, since the model receives additional information about the atmosphere's rate of change and only has to be called half as often.

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.

			$Z_{500}$		$T_{2m}$			$T_{850}$		
	Lead time	W21	HPX64	IFS	W21	HPX64	IFS	W21	HPX64	IFS
ISE	3 days	36.26	21.88	14.91	1.17	0.82	1.02	1.95	1.49	1.35
I RV	5 days	59.01	_41.91	31.30	1.67	1.27	1.27	2.83	2.28	1.96
Ŋ	3 days	0.90	0.96	0.98	0.84	0.92	0.91	0.84	0.91	0.94
ΑC	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 WDCC to our HPX64 model with  $\Delta_t = 3$  h. 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.

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

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



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.

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

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

420

#### 4.2 Eliminating the Need for Boundary-Layer Parameterizations

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

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

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

450

#### 4.3 Iterative Rollouts Over Subseasonal to Annual Time Scales

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

To investigate the stability and drift in model simulations over a full annual cycle, 457 we initialize it using ERA5 data for 00 UTC on 1 June 2017 (together with the 21 UTC 458 fields on 31 May). Using 6 h time steps (with 3 h time resolution), we perform 1460 it-450 erations to generate a 365-day simulation. The three-day running mean of  $Z_{500}$ , aver-460 aged around each latitude, is plotted as a function of latitude and time in Figure 7, along 461 with the corresponding averages from the ERA5 data. Despite being trained to minimize 462 RMSE over a single day and not enforcing any physical constraints, the DLWP-HPX sim-463 ulation responds to the TOA solar forcing to generate the annual cycle reasonably well. 464

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

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

<sup>&</sup>lt;sup>7</sup> 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).

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

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

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

The version of GraphCast from NVIDIA's Earth2MIP gives reasonable results for 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).

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

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

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

A more quantitative assessment of any tendency of our model to distort the atmospheric state by damping or amplifying mid-latitude perturbations at different wavelengths is provided by the plots of the  $Z_{500}$  power spectral density around 45 °N in Figure 10. These spectra are averaged over 208 biweekly forecasts from the 2017-2018 test set; as such, the initial spectrum in black represents the average state of the atmosphere in the 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 536 spectra for wavelengths  $\lambda$  longer than 500 km (roughly 5 grid intervals), but the power 537 in the shorter waves is amplified. Over the next 36 h, there is a gradual reduction in the 538 amplitude at wavelengths  $\lambda < 1800 \,\mathrm{km}$  to yield a spectrum that is modestly damped 539 over the interval  $380 < \lambda < 1800$  km and amplified at the shortest wavelengths. Sur-540 prisingly, the spectral distribution at two days remains essentially unchanged through 541 at least sub-seasonal forecast lead times of eight weeks, which is consistent with the im-542 pression obtained examining images such as those in Figure 9. There is no gradual am-543 plification or loss of amplitude in the simulated atmospheric systems after the first two 544 days. 545

#### 546 5 Conclusion

We have presented an improved CNN-based DLWP-HPX model that stably fore-547 casts atmospheric evolution over a full one-year cycle using a very limited set of prog-548 nostic variables. The number of actual degrees of freedom characterising predictable at-549 mospheric states at forecast lead times beyond 3–5 days is not known, but is far less than 550 the total number of prognostic variables carried at every grid cell in state-of-the-art NWP 551 models. Here, we have demonstrated that realistic atmospheric simulations can be per-552 formed using just seven prognostic variables above each node on a HEALPix mesh with 553 110 km between the nodes. 554

The HEALPix mesh (Gorski et al., 2005) has been used in astronomy for almost 555 two decades, but has previously seen very little use in atmospheric science. The mesh 556 covers the sphere with a hierarchical grid of equal-area cells uniformly spaced along cir-557 cles at constant latitudes. A particularly important advantage of the HEALPix mesh for 558 weather forecasting with CNNs is that it is an "east to the right" mesh, i.e., east has the 559 same orientation in every HEALPix cell. Weather systems tend to travel west-to-east 560 in mid- and high-latitudes and both east-to-west (tropical cyclones) or west-to-east (Madden-561 Julian Oscillation, convectively coupled Kelvin waves) in the tropics. The kernel weights 562 in our convolutional stencils can more economically learn this behavior than on our pre-563 vious cubed sphere mesh in which the eastward orientation across the stencil varies with 564 longitude, particularly on the polar faces. Although switching from a cube-sphere mesh 565 with  $64 \times 64$  cells on each of the six faces to a HEALPix mesh with  $32 \times 32$  cells on 566 each of the 12 faces reduces the total number of grid points covering the sphere by half, 567 it improves the  $Z_{500}$  RMSE error by almost one day at a 4-day forecast lead time (Figure 5). 569

Two other significant improvements to our model architecture were obtained by adding 570 recursion via GRUs and by inverting the standard way channel depth is refined at deeper 571 layers in the U-Net. In contrast to the original U-Net architecture Ronneberger et al. (2015), 572 our channel depth halves instead of doubles as the spatial resolution is also halved in each 573 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 compar-575 atively fewer parameters to describe the simpler set of global weather patterns. Although 576 this modification pushes the U-Net toward the basic ResNet architecture (He et al., 2016), 577 we find the deeper U-Net layers continue to provide significant skill to the forecasts. 578

Additional modest improvements were implemented by switching to the GELU ac-579 tivation function and to  $2 \times 2$  transposed strided convolutions when up-sampling; by in-580 creasing the total number of trainable parameters from  $2.7\,\mathrm{M}$  to  $9.8\,\mathrm{M}$ , adding the  $Z_{250}$ 581 field, increasing the resolution to HPX64, and increasing the time resolution to 3 h (which 582 gives us a 6 h time step). The benefits of 3-h time resolution were only realized when the 583 model included the GRUs. The 3-h time resolution gives a good forecast of the daily cy-584 cle of surface temperature, and the model also learns the difference in the range of that 585 cycle between regions of tropical forest and desert without geo-specific input data. 586

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

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

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

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

#### 623 Appendix A Deep Learning on the HEALPix

#### 624

# A1 Seamless Evolution of Location Invariant Kernels

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

To contrast and emphasize the difficulty that CNN kernels are facing on the cubed sphere mesh, we plot the lines of constant latitude on the six faces of the cubed sphere and on the twelve faces of the HEALPix in Figure A1. Except for the equator, all lines of constant latitude are bent on the cubed sphere, imposing challenges for a limited set of convolution kernels that must evolve location invariant pattern detectors and functions. 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.

propagate eastward motions at the top-left versus the top right corners of the cubed spherefaces.

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

642

#### A2 Technical Implementation Details

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

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

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

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

resolution of the HEALPix mesh and the configuration of the applied convolution (kernel size and dilation). Otherwise, a hierarchy of padding operations would be required to be implemented and considered.

#### 672 Appendix B DLWP Model Details

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

To generate 1-year rollouts for Pangu-Weather, GraphCast, and FourCastNet2 (SFNO), 686 as plotted in Figure 8, we considered the respective public repositories with the pretrained 687 model weights. More concretely, we generated the SFNO Earth2MIP (fcnv2\_sm) and 688 GraphCast Earth2MIP (graphcast) forecasts with NVIDIA's earth2mip package,<sup>8</sup> specif-689 ically developing a custom script for long rollouts.<sup>9</sup> The SFNO Makani forecast, which 690 responds reasonably to the solar forcing by receiving an additional  $\cos \phi$  input (where  $\phi$  is the solar zenith angle), was generated with NVIDIA's Makani package.<sup>10</sup> Interest-692 ingly, the original GraphCast DeepMind code base<sup>11</sup> produced slightly different results 693 and saturated even faster than the Earth2MIP version, which might result from differ-694 ent random seeds. For the DeepMind version of GraphCast, we downloaded the model weights<sup>12</sup> provided through their repository. Pangu-Weather forecasts in 24 h and 3 h res-696 olution (with respective checkpoint files for the  $24 h^{13}$  and  $3 h^{14}$  models) were generated 697 by using the original repositorv.<sup>15</sup> 698

#### **Open Research Section**

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

<sup>&</sup>lt;sup>8</sup> https://github.com/NVIDIA/earth2mip

<sup>&</sup>lt;sup>9</sup> https://github.com/NVIDIA/earth2mip/blob/main/examples/utils/workflows/1\_year\_run.py

<sup>10</sup> https://github.com/NVIDIA/makani

 $<sup>^{11}\,\</sup>tt{https://github.com/google-deepmind/graphcast}$ 

<sup>&</sup>lt;sup>12</sup> https://storage.googleapis.com/dm\_graphcast/params/GraphCast%20-%20ERA5%201979-2017%20 -%20resolution%200.25%20-%20pressure%20levels%2037%20-%20mesh%202to6%20-%20precipitation%

<sup>20</sup>input%20and%20output.npz

<sup>&</sup>lt;sup>13</sup> https://drive.google.com/file/d/11weQlxcn9fG0zKNW8ne1Khr9ehRTI6HP/view

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

<sup>&</sup>lt;sup>15</sup> 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.

							Parameter count		
Layer	$c_{in}$	k	s	d	$\operatorname{RF}$	Output shape	Weights	Biases	Σ
ConvNeXt									
Conv2d	18	1	1	1	$1 \times 1$	(12, 64, 64, 136)	2448	136	2584
Conv2d	18	3	1	1	$3 \times 3$	(12, 64, 64, 544)	88128	544	88672
Conv2d	544	3	1	1	$5 \times 5$	(12, 64, 64, 544)	2663424	544	2663968
Conv2d (1)	544	1	1	1	$5 \times 5$	(12, 64, 64, 136)	73984	136	74120
AvgPool2d	136	2	2		$6 \times 6$	(12, 32, 32, 136)	0	0	0
ConvNeXt									
Conv2d	136	1	1	1	$6 \times 6$	(12, 32, 32, 68)	9248	68	9316
Conv2d	136	3	1	2	$14 \times 14$	(12, 32, 32, 272)	332928	272	333200
Conv2d	272	3	1	<b>2</b>	$22 \times 22$	(12, 32, 32, 272)	665856	272	666128
Conv2d (2)	272	1	1	1	$22 \times 22$	(12, 32, 32, 68)	18496	68	18564
AvgPool2d	68	2	2		$24 \times 24$	(12, 16, 16, 68)	0	0	0
ConvNeXt									
Conv2d	68	1	1	1	$24 \times 24$	(12, 16, 16, 34)	2312	34	2346
Conv2d	68	3	1	4	$56 \times 56$	(12, 16, 16, 136)	83232	136	83368
Conv2d	136	3	1	4	$88 \times 88$	(12, 16, 16, 136)	166464	136	166600
Conv2d	136	1	1	1	$88 \times 88$	(12, 16, 16, 34)	4624	34	4658
ConvNeXt									
Conv2d	34	1	1	1	$88 \times 88$	(12, 16, 16, 68)	2312	68	2380
Conv2d	34	3	1	4	$120 \times 120$	(12, 16, 16, 136)	41616	136	41752
Conv2d	136	3	1	4	$152 \times 152$	(12, 16, 16, 136)	166464	136	166600
Conv2d	136	1	1	1	$152 \times 152$	(12, 16, 16, 68)	9248	68	9316
GRU									
Conv2d	136	1	1	1	$152 \times 152$	(12, 16, 16, 136)	18496	136	18632
Conv2d	136	1	1	1	$152 \times 152$	(12, 16, 16, 68)	9248	68	9316
ConvTrans2d	68	2	2	1	$154 \times 154$	(12, 32, 32, 68)	18496	68	18476
Concat(2)	_					(12, 32, 32, 136)	0	0	0
ConvNeXt									
Conv2d	136	3	1	<b>2</b>	$154 \times 154$	(12, 32, 32, 272)	332928	272	333200
Conv2d	272	3	1	2	$162\times162$	(12, 32, 32, 272)	665856	272	666128
Conv2d	272	1	1	1	$170 \times 170$	(12, 32, 32, 136)	36992	136	37128
GRU									
Conv2d	272	1	1	1	$170\times170$	(12, 32, 32, 272)	73984	272	74256
Conv2d	136	1	1	1	$170\times170$	(12, 32, 32, 136)	36992	136	37128
ConvTrans2d	136	2	2	1	$171 \times 171$	(12, 64, 64, 136)	73984	136	74120
$\texttt{Concat}\ (1)$	_	_	_			(12, 64, 64, 272)	0	0	0
ConvNeXt									
Conv2d	272	1	1	1	$171\times171$	(12, 64, 64, 136)	36992	136	37128
Conv2d	272	3	1	1	$173\times173$	(12, 64, 64, 544)	1331712	544	1332256
Conv2d	544	3	1	1	$175\times175$	(12, 64, 64, 544)	2663424	544	2663968
Conv2d	544	1	1	1	$175\times175$	(12, 64, 64, 136)	73984	136	74120
GRU									
Conv2d	272	1	1	1	$175\times175$	(12, 64, 64, 272)	73984	272	74256
Conv2d	272	1	1	1	$175\times175$	(12, 64, 64, 136)	36992	136	37128
Conv2d	136	1	1	1	$175\times175$	(12, 64, 64, 14)	1904	14	1918
							9816752	6.066	9822818

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

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.

#### 713 Acknowledgments

We would like to thank Mauro Bisson from NVIDIA Corp. for providing optimized CUDA 714 kernels for the HEALPix padding implementation, and Jonathan Weyn who previously 715 implemented a code base on which this work was built. We thank Peter Düben and a 716 second anonymous reviewer for encouraging us to generate and compare the 1-year roll-717 outs for other state-of-the-art DLWP methods and for other valuable suggestions. This work received funding from Deutsche Forschungsgemeinschaft (DFG, German Research 719 Foundation) under Germany's Excellence Strategy EXC 2064 – 390727645 and from the 720 Office of Naval Research under grants N0014-21-1-2827 and N00014-22-1-2807. We thank 721 the Deutscher Akademischer Austauschdienst (DAAD, German Academic Exchange Ser-722 vice) as well as the International Max Planck Research School for Intelligent Systems 723 (IMPRS-IS) for supporting Matthias Karlbauer. Nathaniel was supported by a National 724 Defense Science and Engineering Graduate Fellowship. We are grateful to NVIDIA and 725 Stan Posey for the donation of A100 GPU cards. This research was additionally supported 726 by a grant from the NVIDIA Applied Research Accelerator Program and utilized an NVIDIA 727 DGX-100 Workstation. Moreover, this work benefited substantially from the barrier-free 728 high quality ERA5 dataset provided by the ECMWF. 729

#### 730 Author Roles

Matthias implemented model, training and evaluation routines in PyTorch, as well 731 as the HEALPix-related projection scripts under consideration of the healpy package, 732 and drafted the manuscript together with Dale who supervised this project closely and 733 who also made the model schematic in Figure 2. Nathaniel was involved in discussions 734 about model evolution and code structures and generated Figure 6, Figure 7, and Fig-735 ure 10. Raul was involved in model discussions and generated Figure 9. Thorsten helped 736 with implementing the distributed PyTorch pipeline for multi-GPU training and with 737 accelerating the process pipeline. Noah Brenowitz and Boris Bonev generated the 365days rollouts with the Earth2MIP and Makani packages for SFNO and GraphCast. Mar-739 tin co-supervised this project and helped with proofreading and writing. 740

#### 741 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.
- 745Bauer, P., Dueben, P., Chantry, M., Doblas-Reyes, F., Hoefler, T., McGovern, A.,746& Stevens, B. (2023). Deep learning and a changing economy in weather747and climate prediction. Nature Reviews Earth & Environment, 4(8), 507-748509. Retrieved from https://doi.org/10.1038/s43017-023-00468-z74910.1038/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). Enforc ing 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. <i>Nature</i> . doi: doi.org/
760	10.1038/s41586-023-06185-3
761	Boney, 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 B & Neumann J V (1950) Numerical Integration of the
765	Barotropic Vorticity Equation. <i>Tellus A</i> , 2(4).
766	Chen K Han T Gong J Bai L Ling F Luo J-J Ouvang W (2023)
767	Fengwu: Pushing the skillful global medium-range weather forecast beyond 10
768	days lead. arXiv preprint arXiv:2304.02948.
760	Cho K van Merrienboer B. Gulcebre C. Bougares F. Schwenk H. & Bengio V.
709	(2014) Learning phrase representations using run encoder-decoder for statistical
771	machine translation. In <i>Conference on empirical methods in natural language</i>
772	processing (emplo 2014).
773	Dosovitskiv A Bever 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 ensem-
800	ble 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.

- Krachmalnicoff, N., & Tomasi, M. (2019).Convolutional neural networks on 812 the healpix sphere: a pixel-based algorithm and its application to cmb data 813 analysis. Astronomy & Astrophysics, 628, A129. 814
- Kurth, T., Subramanian, S., Harrington, P., Pathak, J., Mardani, M., Hall, D., ... 815 Anandkumar, A. (2022).Fourcastnet: Accelerating global high-resolution 816 weather forecasting using adaptive fourier neural operators. arXiv preprint 817 arXiv:2208.05419. 818
- Lam, R., Sanchez-Gonzalez, A., Willson, M., Wirnsberger, P., Fortunato, M., Pritzel, 819 A., ... others (2022). Graphcast: Learning skillful medium-range global weather 820 forecasting. arXiv preprint arXiv:2212.12794. 821
- Li, Z., Kovachki, N., Azizzadenesheli, K., Liu, B., Bhattacharya, K., Stuart, A., 822 & Anandkumar, A. (2020).Fourier neural operator for parametric partial 823 differential equations. arXiv preprint arXiv:2010.08895. 824
- Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., ... Guo, B. (2021). Swin trans-825 former: Hierarchical vision transformer using shifted windows. In *Proceedings of* 826 the ieee/cvf international conference on computer vision (pp. 10012-10022). 827
- Liu, Z., Mao, H., Wu, C.-Y., Feichtenhofer, C., Darrell, T., & Xie, S. (2022). A con-828 vnet for the 2020s. In Proceedings of the ieee/cvf conference on computer vision 829 and pattern recognition (pp. 11976–11986). 830
- Lopez-Gomez, I., McGovern, A., Agrawal, S., & Hickey, J. (2022). Global extreme 831 heat forecasting using neural weather models. arXiv preprint arXiv:2205.10972. 832
- Lorenz, E. N. (1969). The predictability of a flow which possesses many scales of mo-833 tion. Tellus, 21(3), 289-307. 834
- Loshchilov, I., & Hutter, F. (2016). Sgdr: Stochastic gradient descent with warm 835 restarts. In International conference on learning representations. 836
- Palmer, T. (2019). The ecmwf ensemble prediction system: Looking back (more than) 837 25 years and projecting forward 25 years. Quarterly Journal of the Royal Mete-838 orological Society, 145, 12–24. 839
- Pathak, J., Subramanian, S., Harrington, P., Raja, S., Chattopadhyay, A., Mardani, 840 M., ... others (2022).Fourcastnet: A global data-driven high-resolution 841 weather model using adaptive fourier neural operators. arXiv preprint 842 arXiv:2202.11214. 843
- Perraudin, N., Defferrard, M., Kacprzak, T., & Sgier, R. (2019). Deepsphere: Efficient spherical convolutional neural network with healpix sampling for cosmological 845 applications. Astronomy and Computing, 27, 130–146. 846

844

847

848

- Pfaff, T., Fortunato, M., Sanchez-Gonzalez, A., & Battaglia, P. W. (2020). Learning mesh-based simulation with graph networks. arXiv preprint arXiv:2010.03409.
- Rasp, S., & Thuerey, N. (2021). Data-driven medium-range weather prediction with a 849 resnet pretrained on climate simulations: A new model for weatherbench. Journal of Advances in Modeling Earth Systems, 13(2), e2020MS002405. 851
- Ronneberger, O., Fischer, P., & Brox, T. (2015). U-net: Convolutional networks for 852 biomedical image segmentation. In International conference on medical image 853 computing and computer-assisted intervention (pp. 234–241). 854
- Scarselli, F., Gori, M., Tsoi, A. C., Hagenbuchner, M., & Monfardini, G. (2008). The 855 graph neural network model. *IEEE transactions on neural networks*, 20(1), 61-856 80. 857
- Scher, S., & Messori, G. (2018). Predicting weather forecast uncertainty with machine 858 learning. Quarterly Journal of the Royal Meteorological Society, 144 (717), 2830-859 2841.
- Scher, S., & Messori, G. (2019). Weather and climate forecasting with neural net-861 works: using GCMs with different complexity as study-ground. Geoscientific 862 Model Development, 12, 2797-2809. 863
- Shen, C., Appling, A. P., Gentine, P., Bandai, T., Gupta, H., Tartakovsky, A., ... 864
- Lawson, K. (2023). Differentiable modelling to unify machine learning and 865 physical models for geosciences. Nature Reviews Earth & Environment, 4(8), 866

867	552-567. Retrieved from https://doi.org/10.1038/s43017-023-00450-9
868	doi: 10.1038/s43017-023-00450-9
869	Thuemmel, J., Karlbauer, M., Otte, S., Zarfl, C., Martius, G., Ludwig, N., others
870	(2023). Inductive biases in deep learning models for weather prediction. $arXiv$
871	$preprint \ arXiv: 2304.04664.$
872	Tobler, W. R. (1970). A computer movie simulating urban growth in the detroit re-
873	gion. Economic geography, $46(\sup 1)$ , $234-240$ .
874	Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N.,
875	Polosukhin, I. (2017). Attention is all you need. Advances in neural information
876	$processing \ systems, \ 30.$
877	Vitart, F. (2004). Monthly forecasting at ECMWF. Monthly Weather Review, 132,
878	2761–2779. doi: 10.1175/MWR2826.1
879	Weigel, A. P., Baggenstos, D., Liniger, M. A., Vitart, F., & Appenzeller, C. (2008).
880	Probabilistic Verification of Monthly Temperature Forecasts. Monthly Weather
881	Review, $136$ , $5162-5182$ . doi: $10.1175/2008MWR2551.1$
882	Weyn, J. A., Durran, D. R., & Caruana, R. (2019). Can machines learn to predict
883	weather? using deep learning to predict gridded 500-hpa geopotential height
884	from historical weather data. Journal of Advances in Modeling Earth Systems,
885	11(8), 2680-2693.
886	Weyn, J. A., Durran, D. R., & Caruana, R. (2020). Improving data-driven global
887	weather prediction using deep convolutional neural networks on a cubed sphere.
888	Journal of Advances in Modeling Earth Systems, 12(9), e2020MS002109.
889	Weyn, J. A., Durran, D. R., Caruana, R., & Cresswell-Clay, N. (2021). Sub-seasonal
890	forecasting with a large ensemble of deep-learning weather prediction models.
891	Journal of Advances in Modeling Earth Systems, 13(7), e2021MS002502.
892	Zhou, Z., Rahman Siddiquee, M. M., Tajbakhsh, N., & Liang, J. (2018). Unet++:
893	A nested u-net architecture for medical image segmentation. In <i>Deep learning</i>
894	in medical image analysis and multimodal learning for clinical decision support
895	(pp. 3–11). Springer.