



Cite as: R. Lam *et al.*, *Science*  
10.1126/science.adi2336 (2023).

# Learning skillful medium-range global weather forecasting

Remi Lam<sup>1\*†</sup>, Alvaro Sanchez-Gonzalez<sup>1\*†</sup>, Matthew Willson<sup>1\*†</sup>, Peter Wirnsberger<sup>1†</sup>,  
Meire Fortunato<sup>1†</sup>, Ferran Alet<sup>1†</sup>, Suman Ravuri<sup>1†</sup>, Timo Ewalds<sup>1</sup>, Zach Eaton-Rosen<sup>1</sup>, Weihua Hu<sup>1</sup>,  
Alexander Merose<sup>2</sup>, Stephan Hoyer<sup>2</sup>, George Holland<sup>1</sup>, Oriol Vinyals<sup>1</sup>, Jacklynn Stott<sup>1</sup>, Alexander Pritzel<sup>1</sup>,  
Shakir Mohamed<sup>1\*</sup>, Peter Battaglia<sup>1\*</sup>

<sup>1</sup>Google DeepMind, London, UK. <sup>2</sup>Google Research, Mountain View, CA, USA.

\*Corresponding author. Email: remilam@google.com (R.L.); alvarosg@google.com (A.S.); matthjw@google.com (M.W.); shakir@google.com (S.M.); peterbattaglia@google.com (P.B.)

†These authors contributed equally to this work.

**Global medium-range weather forecasting is critical to decision-making across many social and economic domains. Traditional numerical weather prediction uses increased compute resources to improve forecast accuracy, but does not directly use historical weather data to improve the underlying model. Here, we introduce “GraphCast,” a machine learning-based method trained directly from reanalysis data. It predicts hundreds of weather variables, over 10 days at 0.25° resolution globally, in under one minute. GraphCast significantly outperforms the most accurate operational deterministic systems on 90% of 1380 verification targets, and its forecasts support better severe event prediction, including tropical cyclones tracking, atmospheric rivers, and extreme temperatures. GraphCast is a key advance in accurate and efficient weather forecasting, and helps realize the promise of machine learning for modeling complex dynamical systems.**

It is 05:45 UTC in mid-October, 2022, in Bologna, Italy, at the European Centre for Medium-Range Weather Forecasts (ECMWF)’s new High-Performance Computing Facility, which recently opened for operation. For the past several hours the Integrated Forecasting System (IFS) has been running sophisticated calculations to forecast Earth’s weather over the next days and weeks, and its first predictions have just begun to be disseminated to users. This process repeats every six hours, every day, to supply the world with the most accurate weather forecasts available.

The IFS, and modern weather forecasting more generally, are triumphs of science and engineering. The dynamics of weather systems are among the most complex physical phenomena on Earth, and each day, countless decisions made by individuals, industries, and policymakers depend on accurate weather forecasts, from deciding whether to wear a jacket or to flee a dangerous storm. The dominant approach for weather forecasting today is “numerical weather prediction” (NWP), which involves solving the governing equations of weather using supercomputers. The success of NWP lies in the rigorous and ongoing research practices that provide increasingly detailed descriptions of weather phenomena, and how well NWP scales to greater accuracy with greater computational resources (1, 2). As a result, the accuracy of weather forecasts has increased year after year, to the point

where the path of a hurricane can be predicted many days ahead—a possibility that was unthinkable even a few decades ago.

But while traditional NWP scales well with compute, capitalizing on the vast amount of historical weather data to improve accuracy is not straightforward. Rather, NWP methods are improved by highly trained experts innovating better models, algorithms, and approximations, which can be a time-consuming and costly process.

Machine learning-based weather prediction (MLWP) offers an alternative to traditional NWP, where forecast models can be trained from historical data, including observations and analysis data. This has potential to improve forecast accuracy by capturing patterns in the data which are not easily represented in explicit equations. MLWP also offers opportunities for greater efficiency by exploiting modern deep learning hardware, rather than supercomputers, and striking more favorable speed-accuracy trade-offs. Recently MLWP has helped improve on NWP-based forecasting in regimes where traditional NWP is relatively weak, for example sub-seasonal heat wave prediction (3) and precipitation nowcasting from radar images (4–7), where accurate equations and robust numerical methods are not as available.

In medium-range weather forecasting, i.e., predicting atmospheric variables up to 10 days ahead, NWP-based systems

Downloaded from https://www.science.org on December 20, 2023

like the IFS are still most accurate. The top deterministic operational system in the world is ECMWF's High RESolution forecast (HRES), a configuration of IFS which produces global 10-day forecasts at  $0.1^\circ$  latitude/longitude resolution, in around an hour (8). However, over the past several years, MLWP methods for medium-range forecasting trained on re-analysis data have been steadily advancing, facilitated by benchmarks such as WeatherBench (8). Deep learning architectures based on convolutional neural networks (9–11) and Transformers (12) have shown promising results at latitude/longitude resolutions coarser than  $1.0^\circ$ , and recent works—which use graph neural networks (GNN), Fourier neural operators, and Transformers (13–16)—have reported performance that begins to rival IFS's at  $1.0^\circ$  and  $0.25^\circ$  for a handful of variables, and lead times up to seven days.

## GraphCast

Here we introduce an MLWP approach for global medium-range weather forecasting called “GraphCast,” which produces an accurate 10-day forecast in under a minute on a single Google Cloud TPU v4 device, and supports applications including predicting tropical cyclone tracks, atmospheric rivers, and extreme temperatures.

GraphCast takes as input the two most recent states of Earth's weather—the current time and six hours earlier—and predicts the next state of the weather six hours ahead. A single weather state is represented by a  $0.25^\circ$  latitude/longitude grid ( $721 \times 1440$ ), which corresponds to roughly  $28 \times 28$  km resolution at the equator (Fig. 1A), where each grid point represents a set of surface and atmospheric variables (listed in Table 1). Like traditional NWP systems, GraphCast is autoregressive: it can be “rolled out” by feeding its own predictions back in as input, to generate an arbitrarily long trajectory of weather states (Fig. 1, B and C).

GraphCast is implemented as a neural network architecture, based on GNNs in an “encode-process-decode” configuration (13, 17), with a total of 36.7 million parameters (code, weights and demos can be found at <https://github.com/deepmind/graphcast>). Previous GNN-based learned simulators (18–20) have been very effective at learning the complex dynamics of fluid and other systems modeled by partial differential equations, which supports their suitability for modeling weather dynamics.

The encoder (Fig. 1D) uses a single GNN layer to map variables (normalized to zero-mean unit-variance) represented as node attributes on the input grid to learned node attributes on an internal “multi-mesh” representation.

The multi-mesh (Fig. 1G) is a graph which is spatially homogeneous, with high spatial resolution over the globe. It is defined by refining a regular icosahedron (12 nodes, 20 faces, 30 edges) iteratively six times, where each refinement divides each triangle into four smaller ones (leading to four times

more faces and edges), and reprojecting the nodes onto the sphere. The multi-mesh contains the 40,962 nodes from the highest resolution mesh (which is roughly  $1/25$  the number of latitude/longitude grid points at  $0.25^\circ$ ), and the union of all the edges created in the intermediate graphs, forming a flat hierarchy of edges with varying lengths.

The processor (Fig. 1E) uses 16 unshared GNN layers to perform learned message-passing on the multi-mesh, enabling efficient local and long-range information propagation with few message-passing steps.

The decoder (Fig. 1F) maps the final processor layers learned features from the multi-mesh representation back to the latitude-longitude grid. It uses a single GNN layer, and predicts the output as a residual update to the most recent input state (with output normalization to achieve unit-variance on the target residual). See supplementary materials section 3 for further architectural details.

During model development, we used 39 years (1979–2017) of historical data from ECMWF's ERA5 (21) reanalysis archive. As a training objective, we averaged the mean squared error (MSE) between GraphCast's predicted states over  $N$  autoregressive steps and the corresponding ERA5 states, with the error weighted by vertical level (see supplementary materials eq. 19). The value of  $N$  was increased incrementally from 1 to 12 (i.e., six hours to three days) over the course of training and the gradient of the loss was computed by backpropagation-through-time (22). GraphCast was trained to minimize the training objective using gradient descent which took roughly four weeks on 32 Cloud TPU v4 devices using batch parallelism. See supplementary materials section 4 for further training details.

Consistent with real deployment scenarios, where future information is not available for model development, we evaluated GraphCast on the held out data from the years 2018 onward (see supplementary materials section 5.1).

## Verification methods

We verify GraphCast's forecast skill comprehensively by comparing its accuracy to HRES's on a large number of variables, levels, and lead times. We quantify the respective skills of GraphCast, HRES, and ML baselines with two skill metrics: the root mean square error (RMSE) and the anomaly correlation coefficient (ACC).

Of the 227 variable and level combinations predicted by GraphCast at each grid point, we evaluated its skill versus HRES on 69 of them, corresponding to the 13 levels of WeatherBench (8) and variables (23) from the ECMWF Scorecard (24); see boldface variables and levels in Table 1 and supplementary materials section 1.2 for which HRES cycle was operational during the evaluation period. In addition to the aggregate performance reported in the main text, supplementary materials section 7 provides further detailed evaluations,

including other variables, precipitation, regional performance, latitude and pressure level effects, spectral properties, blurring, biases, comparisons to other ML-based forecasts, and effects of model design choices.

In making these comparisons, two key choices underlie how skill is established: (i) the selection of the ground truth for comparison, and (ii) a careful accounting of the data assimilation windows used to infer this data from observations. We use ERA5 as the ground truth for evaluating GraphCast, since it was trained to take ERA5 data as input and predict ERA5 data as outputs. However, evaluating HRES forecasts against ERA5 would result in nonzero error on the initial forecast step. Instead, we constructed an “HRES forecast at step 0” (HRES-fc0) dataset to use as ground truth for HRES. HRES-fc0 contains the inputs to HRES forecasts at future initializations (see supplementary materials section 1.2), ensuring that each data point is grounded by recent observations and that the zeroth step of HRES forecasts will have zero error.

For a fair comparison, we must ensure that the ERA5 initial conditions for GraphCast were derived from assimilation windows which look no further into the future than those used by HRES. HRES initializations (00z/06z/12z/18z, where 00z means 00:00 UTC in Zulu convention) always assimilate observations 3 hours into the future while ERA5 initializations assimilate observations 9 hours into the future at 00z/12z and 3 hours into the future at 06z/18z. This constrained the choice of initialization times for GraphCast to 06z/18z in all our results. We use the same initializations for HRES when comparing performance up to 3.75 days. Beyond that, HRES archived forecasts are only available from 00z/12z initializations. The transition from 06z/18z to 00z/12z initializations for HRES induces a small discontinuity in our plots that is indicated by a vertically dashed line at the appropriate lead time. Supplementary materials section 5 contains further verification details, including details of the comparisons protocol between GraphCast and HRES (supplementary materials section 5.2), and the effect of initialization lookahead on both models’ performance (supplementary materials section 5.2.2).

## Forecast verification results

We find that GraphCast has greater weather forecasting skill than HRES when evaluated on 10-day forecasts at a horizontal resolution of  $0.25^\circ$  for latitude/longitude and at 13 vertical levels.

Figure 2, A to C, shows how GraphCast (blue lines) outperforms HRES (black lines) on the Z500 (geopotential at 500 hPa) “headline” field in terms of RMSE skill, RMSE skill score (i.e., the normalized RMSE difference between model *A* and baseline *B* defined as  $(RMSE_A - RMSE_B)/(RMSE_B)$ , and ACC skill. Using Z500, which encodes the synoptic-scale pressure

distribution, is common in the literature, as it has strong meteorological importance (8). The plots show GraphCast has better skill scores across all lead times, with a skill score improvement around 7%–14%. Plots for additional headline variables are in supplementary materials section 7.1.

Figure 2D summarizes the RMSE skill scores for all 1380 evaluated variables and pressure levels, across the 10-day forecasts, in a format analogous to the ECMWF Scorecard. The cell colors are proportional to the skill score, where blue indicates GraphCast had better skill and red indicates HRES had higher skill. GraphCast outperformed HRES on 90.3% of the 1380 targets, and significantly ( $p \leq 0.05$ , nominal sample size  $n \in \{729, 730\}$ ) outperformed HRES on 89.9% of targets. See supplementary materials section 5.4 for methodology and table S4 for *p*-values, test statistics and effective sample sizes.

The regions of the atmosphere in which HRES had better performance than GraphCast (top rows in red in the scorecards), were disproportionately localized in the stratosphere, and had the lowest training loss weight (see supplementary materials section 7.2.2). When excluding the 50 hPa level, GraphCast significantly outperforms HRES on 96.9% of the remaining 1280 targets. When excluding levels 50 and 100 hPa, GraphCast significantly outperforms HRES on 99.7% of the 1180 remaining targets. When conducting per region evaluations, we found the previous results to generally hold across the globe, as detailed in figs. S14 to S16.

We found that increasing the number of autoregressive steps in the MSE loss improves GraphCast performance at longer lead time (see supplementary materials section 7.3.2). It also encourages GraphCast to blur to a degree at longer lead times (see fig. S38), which means its forecasts will lie somewhere in between a traditional deterministic forecast, and an ensemble mean. HRES’s underlying physical equations, however, do not lead to blurred predictions. To assess whether GraphCast’s relative advantage over HRES on RMSE skill is due to blurrier forecasts better optimizing RMSE, we artificially blurred HRES’s forecasts with blurring filters. We fit filters for GraphCast and HRES by minimizing the RMSE between filtered predictions and the models’ respective ground truths. We found that RMSE-optimized blurring applied to GraphCast has greater skill than analogous blurring applied to HRES on 88.0% of our 1380 verification targets, which is generally consistent with our above conclusions (see supplementary materials section 7.4). Still, blurrier forecasts may not be desirable for some applications, which we discuss further in the Conclusions section.

We also compared GraphCast’s performance to the top competing ML-based weather model, Pangu-Weather (16), and found GraphCast outperformed it on 99.2% of the 252 targets they presented (see supplementary materials section 6 for details).



### Severe event forecasting results

Beyond evaluating GraphCast's forecast skill against HRES's on a wide range of variables and lead times, we also evaluate how its forecasts support predicting severe events, including tropical cyclones tracks, atmospheric rivers, and extreme temperature. These are key downstream applications for which GraphCast is not specifically trained, but which are very important for human activity.

### Tropical cyclone tracks

Improving the accuracy of tropical cyclone tracking can help avoid injury and loss of life, as well as reducing economic harm (25). A cyclone's existence and trajectory is predicted by applying a tracking algorithm to forecasts of geopotential ( $Z$ ), horizontal wind ( $10U/10V$ ,  $U/V$ ), and mean sea-level pressure ( $MSL$ ). We implemented a tracking algorithm based on ECMWF's published protocols (26) and applied it to GraphCast's forecasts, to produce cyclone track predictions (see supplementary materials section 8.1). As a baseline for comparison, we used the operational tracks obtained from HRES's  $0.1^\circ$  forecasts with ECMWF's own tracker, stored in the TIGGE archive (27, 28). Each model using the tracker leading to its best performance, we measured errors for both models against the tracks from IBTrACS (29, 30), a separate reanalysis dataset of cyclone tracks aggregated from various analysis and observational sources. Consistent with established evaluation of tropical cyclone prediction (26), we evaluate all tracks when both GraphCast and HRES detect a cyclone, ensuring that both models are evaluated on the same events, and verify that each model's true-positive rates are similar.

Figure 3A shows GraphCast has lower median track error than HRES over 2018–2021 (median was chosen to resist outliers). As per-track errors for HRES and GraphCast are correlated, we also measured the per-track paired error difference between the two models and found that GraphCast is significantly better than HRES for lead time 18 hours to 4.75 days, as shown in Fig. 3B. The error bars show the bootstrapped 95% confidence intervals for the median (see supplementary materials section 8.1 for details).

### Atmospheric rivers

Atmospheric rivers are narrow regions of the atmosphere which are responsible for the majority of the poleward water vapor transport across the mid-latitudes and generate 30%–65% of annual precipitation on the U.S. West Coast (31). Their strength can be characterized by the vertically integrated water vapor transport  $IVT$  (32, 33), indicating whether an event will provide beneficial precipitation or be associated with catastrophic damage (34).  $IVT$  can be computed from the non-linear combination of the horizontal wind speed  $U$  and  $V$  and specific humidity ( $Q$ ), which GraphCast predicts. We

evaluate GraphCast forecasts over coastal North America and the Eastern Pacific during cold months (Oct–Apr), when atmospheric rivers are most frequent. Despite not being specifically trained to characterize atmospheric rivers, Fig. 3C shows that GraphCast improves the prediction of  $IVT$  compared to HRES, from 25% at short lead time, to 10% at longer horizons (see supplementary materials section 8.2 for details).

### Extreme heat and cold

Extreme heat and cold are characterized by large anomalies with respect to typical climatology (3, 35, 36), which can be dangerous and disrupt human activities. We evaluate the skill of HRES and GraphCast in predicting events above the top 2% climatology across location, time of day, and month of the year, for  $2T$  at 12-hour, 5-day, and 10-day lead times, for land regions across northern and southern hemisphere over their respective summer months. We plot precision-recall curves (37) to reflect different possible trade-offs between reducing false positives (high precision) and reducing false negatives (high recall). For each forecast, we obtain the curve by varying a “gain” parameter that scales the  $2T$  forecast's deviations with respect to the median climatology.

Figure 3D shows GraphCast's precision-recall curves are above HRES's for 5- and 10-day lead times, suggesting GraphCast's forecasts are generally superior than HRES at extreme classification over longer horizons. By contrast, HRES has better precision-recall at the 12-hour lead time, which is consistent with the  $2T$  skill score of GraphCast over HRES being near zero, as shown in Fig. 2D. We generally find these results to be consistent across other variables relevant to extreme heat, such as  $T850$  and  $Z500$  (36), other extreme thresholds (5%, 2% and 0.5%), and extreme cold forecasting in winter. See supplementary materials section 8.3 for details.

### Effect of training data recency

GraphCast can be re-trained periodically with recent data, which in principle allows it to capture weather patterns that change over time, such as the effects of climate change, and long climate oscillations. We trained four variants of GraphCast, from scratch, with data that always began in 1979, but ended in 2017, 2018, 2019, and 2020, respectively (we label the variant ending in 2017 as “GraphCast: <2018,” etc.). We compared their performances to HRES on 2021 test data.

Figure 4 shows the skill scores (normalized by GraphCast:<2018) of the four variants and HRES, for  $Z500$ . We found that while GraphCast's performance when trained up to before 2018 is still competitive with HRES in 2021, training it up to before 2021 further improves its skill scores (see supplementary materials section 7.1.3). We speculate this recency effect allows recent weather trends to be captured to improve accuracy. This shows that GraphCast's performance

can be improved by re-training on more recent data.

## Conclusions

GraphCast's forecast skill and efficiency compared to HRES shows MLWP methods are now competitive with traditional weather forecasting methods. Additionally, GraphCast's performance on severe event forecasting, which it was not directly trained for, demonstrates its robustness and potential for downstream value. We believe this marks a turning point in weather forecasting, which helps open new avenues to strengthen the breadth of weather-dependent decision-making by individuals and industries, by making cheap prediction more accurate, more accessible, and suitable for specific applications.

With 36.7 million parameters, GraphCast is a relatively small model by modern ML standards, chosen to keep the memory footprint tractable. And while HRES is released on 0.1° resolution, 137 levels, and up to 1 hour time steps, GraphCast operated on 0.25° latitude-longitude resolution, 37 vertical levels, and 6 hour time steps, because of the ERA5 training data's native 0.25° resolution, and engineering challenges in fitting higher resolution data on hardware. Generally GraphCast should be viewed as a family of models, with the current version being the largest we can practically fit under current engineering constraints, but which have potential to scale much further in the future with greater compute resources and higher resolution data.

One key limitation of our approach is in how uncertainty is handled. We focused on deterministic forecasts and compared against HRES, but the other pillar of ECMWF's IFS, the ensemble forecasting system, ENS, is especially important for quantifying the probability of extreme events and as the skill of the forecast decreases at longer lead times. The non-linearity of weather dynamics means there is increasing uncertainty at longer lead times, which is not well-captured by a single deterministic forecast. ENS addresses this by generating multiple, stochastic forecasts, which approximate a predictive distribution over future weather, however generating multiple forecasts is expensive. By contrast, GraphCast's MSE training objective encourages it to spatially blur its predictions in the presence of uncertainty, which may not be desirable for some applications where knowing tail, or joint, probabilities of events is important. Building probabilistic forecasts that model uncertainty more explicitly, along the lines of ensemble forecasts, is a crucial next step.

It is important to emphasize that data-driven MLWP relies critically on large quantities of data and their quality, which in the case of models trained on reanalysis, depends on the fidelity of NWP. Therefore, rich high-quality data sources like ECMWF's MARS archive (38) are invaluable. Our approach should not be regarded as a replacement for traditional weather forecasting methods, which have been

developed for decades, rigorously tested in many real-world contexts, and offer many features we have not yet explored. Rather our work should be interpreted as evidence that MLWP is able to meet the challenges of real-world forecasting problems and has potential to complement and improve the current best methods.

Beyond weather forecasting, GraphCast can open new directions for other important geo-spatiotemporal forecasting problems, including climate and ecology, energy, agriculture, and human and biological activity, as well as other complex dynamical systems. We believe that learned simulators, trained on rich, real-world data, will be crucial in advancing the role of machine learning in the physical sciences.

## REFERENCES AND NOTES

1. S. G. Benjamin, J. M. Brown, G. Brunet, P. Lynch, K. Saito, T. W. Schlatter, 100 years of progress in forecasting and NWP applications. *Meteorol. Monogr.* **59**, 13.1–13.67 (2019). [doi:10.1175/AMSMONOGRAPHSD-18-0020.1](https://doi.org/10.1175/AMSMONOGRAPHSD-18-0020.1)
2. P. Bauer, A. Thorpe, G. Brunet, The quiet revolution of numerical weather prediction. *Nature* **525**, 47–55 (2015). [doi:10.1038/nature14956](https://doi.org/10.1038/nature14956) [Medline](#)
3. I. Lopez-Gomez, A. McGovern, S. Agrawal, J. Hickey, Global extreme heat forecasting using neural weather models. *Artif. Intell. Earth Syst.* **2**, e220035 (2022).
4. X. Shi, Z. Gao, L. Lausen, H. Wang, D.-Y. Yeung, W.-k. Wong, W.-c. Woo, Deep learning for precipitation nowcasting: A benchmark and a new model. *Adv. Neural Inf. Process. Syst.* **30**, 5617–5627 (2017).
5. C. K. Sørderby, L. Espeløt, J. Heek, M. Dehghani, A. Oliver, T. Salimans, S. Agrawal, J. Hickey, N. Kalchbrenner, Metnet: A neural weather model for precipitation forecasting. [arXiv:2003.12140](https://arxiv.org/abs/2003.12140) [cs.LG] (2020).
6. S. Ravuri, K. Lenc, M. Willson, D. Kangin, R. Lam, P. Mirowski, M. Fitzsimons, M. Athanassiadou, S. Kashem, S. Madge, R. Prudden, A. Mandhane, A. Clark, A. Brock, K. Simonyan, R. Hadsell, N. Robinson, E. Clancy, A. Arribas, S. Mohamed, Skilful precipitation nowcasting using deep generative models of radar. *Nature* **597**, 672–677 (2021). [doi:10.1038/s41586-021-03854-z](https://doi.org/10.1038/s41586-021-03854-z) [Medline](#)
7. L. Espeløt, S. Agrawal, C. Sørderby, M. Kumar, J. Heek, C. Bromberg, C. Gazen, R. Carver, M. Andrychowicz, J. Hickey, A. Bell, N. Kalchbrenner, Deep learning for twelve hour precipitation forecasts. *Nat. Commun.* **13**, 5145 (2022). [doi:10.1038/s41467-022-32483-x](https://doi.org/10.1038/s41467-022-32483-x) [Medline](#)
8. S. Rasp, P. D. Dueben, S. Scher, J. A. Weyn, S. Mouatadid, N. Thuerey, WeatherBench: A benchmark data set for data-driven weather forecasting. *J. Adv. Model. Earth Syst.* **12**, e2020MS002203 (2020). [doi:10.1029/2020MS002203](https://doi.org/10.1029/2020MS002203)
9. J. A. Weyn, D. R. Durran, R. Caruana, Can machines learn to predict weather? Using deep learning to predict gridded 500-hPa geopotential height from historical weather data. *J. Adv. Model. Earth Syst.* **11**, 2680–2693 (2019). [doi:10.1029/2019MS001705](https://doi.org/10.1029/2019MS001705)
10. J. A. Weyn, D. R. Durran, R. Caruana, Improving data-driven global weather prediction using deep convolutional neural networks on a cubed sphere. *J. Adv. Model. Earth Syst.* **12**, e2020MS002109 (2020). [doi:10.1029/2020MS002109](https://doi.org/10.1029/2020MS002109)
11. S. Rasp, N. Thuerey, Data-driven medium-range weather prediction with a resnet pretrained on climate simulations: A new model for WeatherBench. *J. Adv. Model. Earth Syst.* **13**, e2020MS002405 (2021). [doi:10.1029/2020MS002405](https://doi.org/10.1029/2020MS002405)
12. T. Nguyen, J. Brandstetter, A. Kapoor, J. K. Gupta, A. Grover, ClimaX: A foundation model for weather and climate. [arXiv:2301.10343](https://arxiv.org/abs/2301.10343) [cs.LG] (2023).
13. R. Keisler, Forecasting global weather with graph neural networks. [arXiv:2202.07575](https://arxiv.org/abs/2202.07575) [physics.ao-ph] (2022).
14. J. Pathak, S. Subramanian, P. Harrington, S. Raja, A. Chattopadhyay, M. Mardani, T. Kurth, D. Hall, Z. Li, K. Azizzadenesheli, P. Hassanzadeh, K. Kashinath, A. Anandkumar, FourCastNet: A global data-driven high-resolution weather model using adaptive fourier neural operators. [arXiv:2202.11214](https://arxiv.org/abs/2202.11214) [physics.ao-ph] (2022).
15. T. Kurth, S. Subramanian, P. Harrington, J. Pathak, M. Mardani, D. Hall, A. Miele, K. Kashinath, A. Anandkumar, FourCastNet: Accelerating global high-resolution weather forecasting using adaptive fourier neural operators. [arXiv:2208.05419](https://arxiv.org/abs/2208.05419)

- [physics.ao-ph] (2022).
16. K. Bi, L. Xie, H. Zhang, X. Chen, X. Gu, Q. Tian, Pangu-Weather: A 3D high-resolution model for fast and accurate global weather forecast. [arXiv:2211.02556](https://arxiv.org/abs/2211.02556) [physics.ao-ph] (2022).
  17. P. W. Battaglia, J. B. Hamrick, V. Bapst, A. Sanchez-Gonzalez, V. Zambaldi, M. Malinowski, A. Tacchetti, D. Raposo, A. Santoro, R. Faulkner, C. Gulcehre, F. Song, A. Ballard, J. Gilmer, G. Dahl, A. Vaswani, K. Allen, C. Nash, V. Langston, C. Dyer, N. Heess, D. Wierstra, P. Kohli, M. Botvinick, O. Vinyals, Y. Li, R. Pascanu, Relational inductive biases, deep learning, and graph networks. [arXiv:1806.01261](https://arxiv.org/abs/1806.01261) [cs.LG] (2018).
  18. A. Sanchez-Gonzalez, J. Godwin, T. Pfaff, R. Ying, J. Leskovec, P. Battaglia, Learning to simulate complex physics with graph networks. *Proc. Mach. Learn. Res.* **119**, 8459–8468 (2020).
  19. T. Pfaff, M. Fortunato, A. Sanchez-Gonzalez, P. Battaglia, “Learning mesh-based simulation with graph networks,” International Conference on Learning Representations (ICLR 2021), 3 to 7 May 2021.
  20. F. Alet, A. K. Jeewajee, M. B. Villalonga, A. Rodriguez, T. Lozano-Perez, L. Kaelbling, Graph element networks: adaptive, structured computation and memory. *Proc. Mach. Learn. Res.* **97**, 212–222 (2019).
  21. H. Hersbach, B. Bell, P. Berrisford, S. Hirahara, A. Horányi, J. Muñoz-Sabater, J. Nicolas, C. Peubey, R. Radu, D. Schepers, A. Simmons, C. Soci, S. Abdalla, X. Abellan, G. Balsamo, P. Bechtold, G. Biavati, J. Bidlot, M. Bonavita, G. De Chiara, P. Dahlgren, D. Dee, M. Diamantakis, R. Dragani, J. Flemming, R. Forbes, M. Fuentes, A. Geer, L. Haimberger, S. Healy, R. J. Hogan, E. Hólm, M. Janisková, S. Keeley, P. Laloyaux, P. Lopez, C. Lupu, G. Radnoti, P. de Rosnay, I. Rozum, F. Vamborg, S. Villaume, J.-N. Thépaut, The ERA5 global reanalysis. *Q. J. R. Meteorol. Soc.* **146**, 1999–2049 (2020). [doi:10.1002/qj.3803](https://doi.org/10.1002/qj.3803)
  22. P. Werbos, Backpropagation through time: What it does and how to do it. *Proc. IEEE* **78**, 1550–1560 (1990). [doi:10.1109/5.58337](https://doi.org/10.1109/5.58337)
  23. Because precipitation in ERA5 has known biases (39), no development decision for GraphCast was made to improve performance on precipitation and GraphCast simply uses precipitation as an auxiliary input/output. Note that precipitation is sparse and non-Gaussian and would have possibly required different modeling decisions than the other variables. Additionally, precipitation is not available in the HRES analysis products in a form amenable to our evaluation protocol (see next paragraphs). Thus, any claim about precipitation prediction is left out of the scope of this work, and we show precipitation evaluation using a different protocol in supplementary materials section 7.1.4 for completeness only.
  24. T. Haiden, M. Janousek, J. Bidlot, R. Buizza, L. Ferranti, F. Prates, F. Vitart, “Evaluation of ECMWF forecasts, including the 2018 upgrade,” ECMWF Technical Memorandum No. 831 (European Centre for Medium-Range Weather Forecasts, 2018); <https://doi.org/10.21957/ldw15ckqj>
  25. A. B. Martinez, Forecast accuracy matters for hurricane damage. *Econometrics* **8**, 18 (2020). [doi:10.3390/econometrics8020018](https://doi.org/10.3390/econometrics8020018)
  26. L. Magnusson, S. Majumdar, R. Emerton, D. Richardson, M. Alonso-Balmaseda, C. Baugh, P. Bechtold, J. Bidlot, A. Bonanni, M. Bonavita, N. Bormann, A. Brown, P. Browne, H. Carr, M. Dahoui, G. De Chiara, M. Diamantakis, D. Duncan, S. English, R. Forbes, A. Geer, T. Haiden, S. Healy, T. Hewson, B. Ingleby, M. Janousek, C. Kuehnlein, S. Lang, S.-J. Lock, T. McNally, K. Mogensen, F. Pappenberger, I. Polichtchouk, F. Prates, C. Prudhomme, F. Rabier, P. de Rosnay, T. Quintino, M. Rennie, H. Tittley, F. Vana, F. Vitart, F. Warrick, N. Wedi, E. Zsoter, “Tropical cyclone activities at ECMWF,” ECMWF Technical Memorandum No. 888 (European Centre for Medium-Range Weather Forecasts, 2021); <https://doi.org/10.21957/zxxzygywv>.
  27. P. Bougeault, Z. Toth, C. Bishop, B. Brown, D. Burridge, D. H. Chen, B. Ebert, M. Fuentes, T. M. Hamill, K. Mylne, J. Nicolau, T. Paccagnella, Y.-Y. Park, D. Parsons, B. Raoult, D. Schuster, P. S. Dias, R. Swinbank, Y. Takeuchi, W. Tennant, L. Wilson, S. Worley, The THORPEX interactive grand global ensemble. *Bull. Am. Meteorol. Soc.* **91**, 1059–1072 (2010). [doi:10.1175/2010BAMS2853.1](https://doi.org/10.1175/2010BAMS2853.1)
  28. R. Swinbank, M. Kyouda, P. Buchanan, L. Froude, T. M. Hamill, T. D. Hewson, J. H. Keller, M. Matsueda, J. Methven, F. Pappenberger, M. Scheuerer, H. A. Tittley, L. Wilson, M. Yamaguchi, The TIGGE project and its achievements. *Bull. Am. Meteorol. Soc.* **97**, 49–67 (2016). [doi:10.1175/BAMS-D-13-00191.1](https://doi.org/10.1175/BAMS-D-13-00191.1)
  29. K. R. Knapp, M. C. Kruk, D. H. Levinson, H. J. Diamond, C. J. Neumann, The international best track archive for climate stewardship (IBTrACS) unifying tropical cyclone data. *Bull. Am. Meteorol. Soc.* **91**, 363–376 (2010). [doi:10.1175/2009BAMS2755.1](https://doi.org/10.1175/2009BAMS2755.1)
  30. K. R. Knapp, H. J. Diamond, J. P. Kossin, M. C. Kruk, C. J. Schreck III, International Best Track Archive for Climate Stewardship (IBTrACS) Project, version 4, NOAA National Centers for Environmental Information (2018); <https://doi.org/10.25921/82ty-9e16>.
  31. W. Chapman, A. Subramanian, L. Delle Monache, S. Xie, F. Ralph, Improving atmospheric river forecasts with machine learning. *Geophys. Res. Lett.* **46**, 10627–10635 (2019). [doi:10.1029/2019GL083662](https://doi.org/10.1029/2019GL083662)
  32. P. J. Neiman, F. M. Ralph, G. A. Wick, J. D. Lundquist, M. D. Dettinger, Meteorological characteristics and overland precipitation impacts of atmospheric rivers affecting the west coast of North America based on eight years of SSM/I satellite observations. *J. Hydrometeorol.* **9**, 22–47 (2008). [doi:10.1175/2007JHM855.1](https://doi.org/10.1175/2007JHM855.1)
  33. B. J. Moore, P. J. Neiman, F. M. Ralph, F. E. Barthold, Physical processes associated with heavy flooding rainfall in Nashville, Tennessee, and vicinity during 1–2 May 2010: The role of an atmospheric river and mesoscale convective systems. *Mon. Weather Rev.* **140**, 358–378 (2012). [doi:10.1175/MWR-D-11-00126.1](https://doi.org/10.1175/MWR-D-11-00126.1)
  34. T. W. Corringham, F. M. Ralph, A. Gershunov, D. R. Cayan, C. A. Talbot, Atmospheric rivers drive flood damages in the western United States. *Sci. Adv.* **5**, eaax4631 (2019). [doi:10.1126/sciadv.aax4631](https://doi.org/10.1126/sciadv.aax4631) [Medline](https://pubmed.ncbi.nlm.nih.gov/33111111/)
  35. L. Magnusson, T. Haiden, D. Richardson, “Verification of extreme weather events: Discrete predictands,” ECMWF Technical Memorandum No. 731 (European Centre for Medium-Range Weather Forecasts, 2014); <https://doi.org/10.21957/liql31n2c>.
  36. L. Magnusson, ECMWF confluence wiki: 202208 - Heatwave - UK (2022); <https://confluence.ecmwf.int/display/FCST/202208++Heatwave++UK>.
  37. T. Saito, M. Rehmsmeier, The precision-recall plot is more informative than the ROC plot when evaluating binary classifiers on imbalanced datasets. *PLOS ONE* **10**, e0118432 (2015). [doi:10.1371/journal.pone.0118432](https://doi.org/10.1371/journal.pone.0118432) [Medline](https://pubmed.ncbi.nlm.nih.gov/26061111/)
  38. C. Maass, E. Cuartero, MARS user documentation (2022); <https://confluence.ecmwf.int/display/UDOC/MARS+user+documentation>.
  39. D. A. Lavers, A. Simmons, F. Vamborg, M. J. Rodwell, An evaluation of ERA5 precipitation for climate monitoring. *Q. J. R. Meteorol. Soc.* **148**, 3152–3165 (2022). [doi:10.1002/qj.4351](https://doi.org/10.1002/qj.4351)
  40. D. Fan, G. Jodin, T. R. Consi, L. Bonfiglio, Y. Ma, L. R. Keyes, G. E. Karniadakis, M. S. Triantafyllou, A robotic Intelligent Towing Tank for learning complex fluid-structure dynamics. *Sci. Robot.* **4**, eaay5063 (2019). [doi:10.1126/scirobotics.aay5063](https://doi.org/10.1126/scirobotics.aay5063) [Medline](https://pubmed.ncbi.nlm.nih.gov/33111111/)
  41. T. Ewalds, Source code for GraphCast, google-deepmind/graphcast: Version 0.1, Zenodo (2023); <https://doi.org/10.5281/zenodo.10058758>
  42. S. Malardel, N. Wedi, W. Deconinck, M. Diamantakis, C. Kuehnlein, G. Mozdzyński, M. Hamrud, P. Smolarkiewicz, “A new grid for the IFS,” ECMWF Newsletter, No. 146, Winter 2015/16, pp. 23–28 (European Centre for Medium-Range Weather Forecasts, 2016); <https://doi.org/10.21957/zwdu9u5j>.
  43. D. H. Levinson, H. J. Diamond, K. R. Knapp, M. C. Kruk, E. J. Gibney, Toward a homogenous global tropical cyclone best-track dataset. *Bull. Am. Meteorol. Soc.* **91**, 377–380 (2010).
  44. M. C. Kruk, K. R. Knapp, D. H. Levinson, A technique for combining global tropical cyclone best track data. *J. Atmos. Ocean. Technol.* **27**, 680–692 (2010). [doi:10.1175/2009JTECHA1267.1](https://doi.org/10.1175/2009JTECHA1267.1)
  45. A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, I. Polosukhin, Attention is all you need. *Adv. Neural Inf. Process. Syst.* **30**, 5998–6008 (2017).
  46. P. Veličković, G. Cucurull, A. Casanova, A. Romero, P. Liò, Y. Bengio, Graph attention networks. [arXiv:1710.10903](https://arxiv.org/abs/1710.10903) [stat.ML] (2017).
  47. M. Fortunato, T. Pfaff, P. Wirsberger, A. Pritzel, P. Battaglia, Multiscale MeshGraphNets. [arXiv:2210.00612](https://arxiv.org/abs/2210.00612) [cs.LG] (2022).
  48. P. Battaglia, R. Pascanu, M. Lai, D. Jimenez Rezende, K. Kavukcuoglu, Interaction networks for learning about objects, relations and physics. *Adv. Neural Inf. Process. Syst.* **29**, 4502–4510 (2016).
  49. K. R. Allen, Y. Rubanova, T. Lopez-Guevara, W. Whitney, A. Sanchez-Gonzalez, P. Battaglia, T. Pfaff, Learning rigid dynamics with face interaction graph networks. [arXiv:2212.03574](https://arxiv.org/abs/2212.03574) [cs.LG] (2022).



50. P. Ramachandran, B. Zoph, Q. V. Le, Searching for activation functions. [arXiv:1710.05941](https://arxiv.org/abs/1710.05941) [cs.NE] (2017).
51. J. L. Ba, J. R. Kiros, G. E. Hinton, Layer normalization. [arXiv:1607.06450](https://arxiv.org/abs/1607.06450) [stat.ML] (2016).
52. B. Bell, H. Hersbach, A. Simmons, P. Berrisford, P. Dahlgren, A. Horányi, J. Muñoz-Sabater, J. Nicolas, R. Radu, D. Schepers, C. Soci, S. Villaume, J.-R. Bidlot, L. Haimberger, J. Woollen, C. Buontempo, J.-N. Thépaut, The ERA5 global reanalysis: Preliminary extension to 1950. *Q. J. R. Meteorol. Soc.* **147**, 4186–4227 (2021). [doi:10.1002/qj.4174](https://doi.org/10.1002/qj.4174)
53. I. Loshchilov, F. Hutter, Decoupled weight decay regularization. [arXiv:1711.05101](https://arxiv.org/abs/1711.05101) [cs.LG] (2017).
54. D. P. Kingma, J. Ba, Adam: A method for stochastic optimization. [arXiv:1412.6980](https://arxiv.org/abs/1412.6980) [cs.LG] (2014).
55. I. Loshchilov, F. Hutter, SGDR: Stochastic Gradient Descent with Warm Restarts. [arXiv:1608.03983](https://arxiv.org/abs/1608.03983) [cs.LG] (2017).
56. T. Chen, B. Xu, C. Zhang, C. Guestrin, Training deep nets with sublinear memory cost. [arXiv:1604.06174](https://arxiv.org/abs/1604.06174) [cs.LG] (2016).
57. J. Bradbury, R. Frostig, P. Hawkins, M. J. Johnson, C. Leary, D. Maclaurin, G. Necula, A. Paszke, J. VanderPlas, S. Wanderman-Milne, Q. Zhang, JAX: composable transformations of Python+NumPy programs (2018); <https://github.com/google/jax>.
58. T. Hennigan, T. Cai, T. Norman, I. Babuschkin, Haiku: Sonnet for JAX (2020); <https://github.com/deepmind/dm-haiku>.
59. J. Godwin, T. Keck, P. Battaglia, V. Bapst, T. Kipf, Y. Li, K. Stachenfeld, P. Veličković, A. Sanchez-Gonzalez, Jraph: A library for graph neural networks in JAX (2020); <https://github.com/deepmind/jraph>.
60. I. Babuschkin, K. Baumli, A. Bell, S. Bhupatiraju, J. Bruce, P. Buchlovsky, D. Budden, T. Cai, A. Clark, I. Danihelka, C. Fantacci, J. Godwin, C. Jones, R. Hemsley, T. Hennigan, M. Hessel, S. Hou, S. Kapturowski, T. Keck, I. Kemaev, M. King, M. Kunesch, L. Martens, H. Merzic, V. Mikulik, T. Norman, J. Quan, G. Papamakarios, R. Ring, F. Ruiz, A. Sanchez, R. Schneider, E. Sezener, S. Spencer, S. Srinivasan, L. Wang, W. Stokowiec, F. Viola, The DeepMind JAX Ecosystem (2020); <https://github.com/deepmind>.
61. S. Hoyer, J. Hamman, xarray: N-D labeled arrays and datasets in Python. *J. Open Res. Softw.* **5**, 10 (2017). [doi:10.5334/jors.148](https://doi.org/10.5334/jors.148)
62. ECMWF, IFS Documentation CY45R1 – Part II: Data assimilation (European Centre for Medium-Range Weather Forecasts, 2018); <https://doi.org/10.21957/a3ri44ig4>.
63. J. R. Driscoll, D. M. Healy, Computing fourier transforms and convolutions on the 2-sphere. *Adv. Appl. Math.* **15**, 202–250 (1994). [doi:10.1006/aama.1994.1008](https://doi.org/10.1006/aama.1994.1008)
64. A. J. Geer, Significance of changes in medium-range forecast scores. *Tellus A Dyn. Meteorol. Oceanogr.* **68**, 30229 (2016). [doi:10.3402/tellusa.v68.30229](https://doi.org/10.3402/tellusa.v68.30229)
65. T. Haiden, M. Janousek, F. Vitart, L. Ferranti, F. Prates, “Evaluation of ECMWF forecasts, including the 2019 upgrade,” ECMWF Technical Memorandum No. 853 (European Centre for Medium-Range Weather Forecasts, 2019); <https://doi.org/10.21957/mlvapkke>.
66. T. Haiden, M. Janousek, F. Vitart, Z. Ben-Bouallegue, L. Ferranti, C. Prates, D. Richardson, “Evaluation of ECMWF forecasts, including the 2020 upgrade,” ECMWF Technical Memorandum No. 880 (European Centre for Medium-Range Weather Forecasts, 2021); <https://doi.org/10.21957/6njp8byz4>.
67. T. Haiden, M. Janousek, F. Vitart, Z. Ben-Bouallegue, L. Ferranti, F. Prates, “Evaluation of ECMWF forecasts, including the 2021 upgrade,” ECMWF Technical Memorandum No. 884 (European Centre for Medium-Range Weather Forecasts, 2021); <https://doi.org/10.21957/90pgicjk4>.
68. T. Haiden, M. Janousek, F. Vitart, Z. Ben-Bouallegue, L. Ferranti, F. Prates, D. Richardson, “Evaluation of ECMWF forecasts, including the 2021 upgrade,” ECMWF Technical Memorandum No. 902 (European Centre for Medium-Range Weather Forecasts, 2022); <https://doi.org/10.21957/xqnu5o3p>.
69. M. J. Rodwell, D. S. Richardson, T. D. Hewson, T. Haiden, A new equitable score suitable for verifying precipitation in numerical weather prediction. *Q. J. R. Meteorol. Soc.* **136**, 1344–1363 (2010). [doi:10.1002/qj.656](https://doi.org/10.1002/qj.656)
70. T. Haiden, M. J. Rodwell, D. S. Richardson, A. Okagaki, T. Robinson, T. Hewson, Intercomparison of global model precipitation forecast skill in 2010/11 using the seeps score. *Mon. Weather Rev.* **140**, 2720–2733 (2012). [doi:10.1175/MWR-D-11-00301.1](https://doi.org/10.1175/MWR-D-11-00301.1)
71. R. North, M. Trueman, M. Mittermaier, M. J. Rodwell, An assessment of the SEEPS and SEDI metrics for the verification of 6 h forecast precipitation accumulations. *Meteorol. Appl.* **20**, 164–175 (2013). [doi:10.1002/met.1405](https://doi.org/10.1002/met.1405)
72. B. D. Santer, R. Sausen, T. M. L. Wigley, J. S. Boyle, K. AchutaRao, C. Doutriaux, J. E. Hansen, G. A. Meehl, E. Roeckner, R. Ruedy, G. Schmidt, K. E. Taylor, Behavior of tropopause height and atmospheric temperature in models, reanalyses, and observations: Decadal changes. *J. Geophys. Res.* **108**, ACL 1-1–ACL 1-22 (2003). [doi:10.1029/2002JD002258](https://doi.org/10.1029/2002JD002258)
73. ECMWF, IFS Documentation CY41R2 – Part III: Dynamics and numerical procedures (European Centre for Medium-Range Weather Forecasts, 2016); <https://doi.org/10.21957/83wouu8Q>.
74. B. Devaraju, “Understanding filtering on the sphere: Experiences from filtering GRACE data,” thesis, University of Stuttgart (2015).
75. H. T. Taylor, B. Ward, M. Willis, W. Zaleski, “The Saffir-Simpson hurricane wind scale” (Atmospheric Administration, Washington, DC, USA, 2010).

## ACKNOWLEDGMENTS

In alphabetical order, we thank Kelsey Allen, Charles Blundell, Matt Botvinick, Zied Ben Bouallegue, Michael Brenner, Rob Carver, Matthew Chantry, Marc Deisenroth, Peter Deuben, Marta Garnelo, Ryan Keisler, Dmitrii Kochkov, Christopher Mattern, Piotr Mirowski, Peter Norgaard, Ilan Price, Chongli Qin, Sébastien Racanière, Stephan Rasp, Yulia Rubanova, Kunal Shah, Jamie Smith, Daniel Worrall, and countless others at Alphabet and ECMWF for advice and feedback on our work. We also thank ECMWF for providing invaluable datasets to the research community. The style of the opening paragraph was inspired by (40). **Funding:** All research in this study was funded by Google DeepMind and Alphabet. There was no external funding. **Author contributions:** Conceptualization: R.L., A.S., M.W., S.M., P.B. Data curation: R.L., A.S., M.W., A.M., P.B. Formal analysis: R.L., A.S., M.W., P.W., M.F., F.A., S.R., T.E., Z.E., W.H., A.P., S.M., P.B. Investigation: R.L., A.S., M.W., P.W., M.F., F.A., S.R., T.E., Z.E., W.H., A.P., S.M., P.B. Methodology: R.L., A.S., M.W., P.W., M.F., F.A., S.R., T.E., Z.E., W.H., A.P., S.M., P.B. Project administration: R.L., A.S., M.W., G.H., O.V., J.S., S.M., P.B. Software: R.L., A.S., M.W., P.W., M.F., F.A., S.R., T.E., Z.E., W.H., A.P., P.B. Supervision: R.L., S.M., P.B. Validation: R.L., A.S., M.W., P.W., M.F., F.A., S.R., T.E., Z.E., W.H., A.P., S.M., P.B. Visualization: R.L., A.S., M.W., F.A., P.B. Writing - original draft: R.L., A.S., M.W., P.W., M.F., F.A., S.R., T.E., Z.E., S.H., A.P., S.M., P.B. Writing - review & editing: R.L., A.S., M.W., P.W., M.F., F.A., S.R., T.E., Z.E., S.H., A.P., S.M., P.B. **Competing interests:** This work was done in the course of employment at Google DeepMind, with no other competing financial interests. A.S., R.L., P.B., M.W., P.W., M.F. and A.P. have filed a provisional patent application relating to machine learning for learned medium-range global weather forecasting (US Provisional App. no. US63/435,163). **Data and materials availability:** GraphCast’s code and trained weights are publicly available on GitHub <https://github.com/deepmind/graphcast> (41). This work used publicly available data from the European Centre for Medium Range Forecasting (ECMWF). We use the ECMWF archive (expired real-time) products for ERA5, HRES and TIGGE products, whose use is governed by the Creative Commons Attribution 4.0 International (CC BY 4.0). We use IBTrACS Version 4 from <https://www.ncei.noaa.gov/products/international-best-track-archive> and reference (29, 30) as required. **License information:** Copyright © 2023 the authors, some rights reserved; exclusive licensee American Association for the Advancement of Science. No claim to original US government works. <https://www.science.org/about/science-licenses-journal-article-reuse>

## SUPPLEMENTARY MATERIALS

[science.org/doi/10.1126/science.adi2336](https://science.org/doi/10.1126/science.adi2336)

Materials and Methods

Supplementary Text

Figs. S1 to S53

Tables S1 to S4

References (42–75)

Submitted 18 April 2023; accepted 1 November 2023

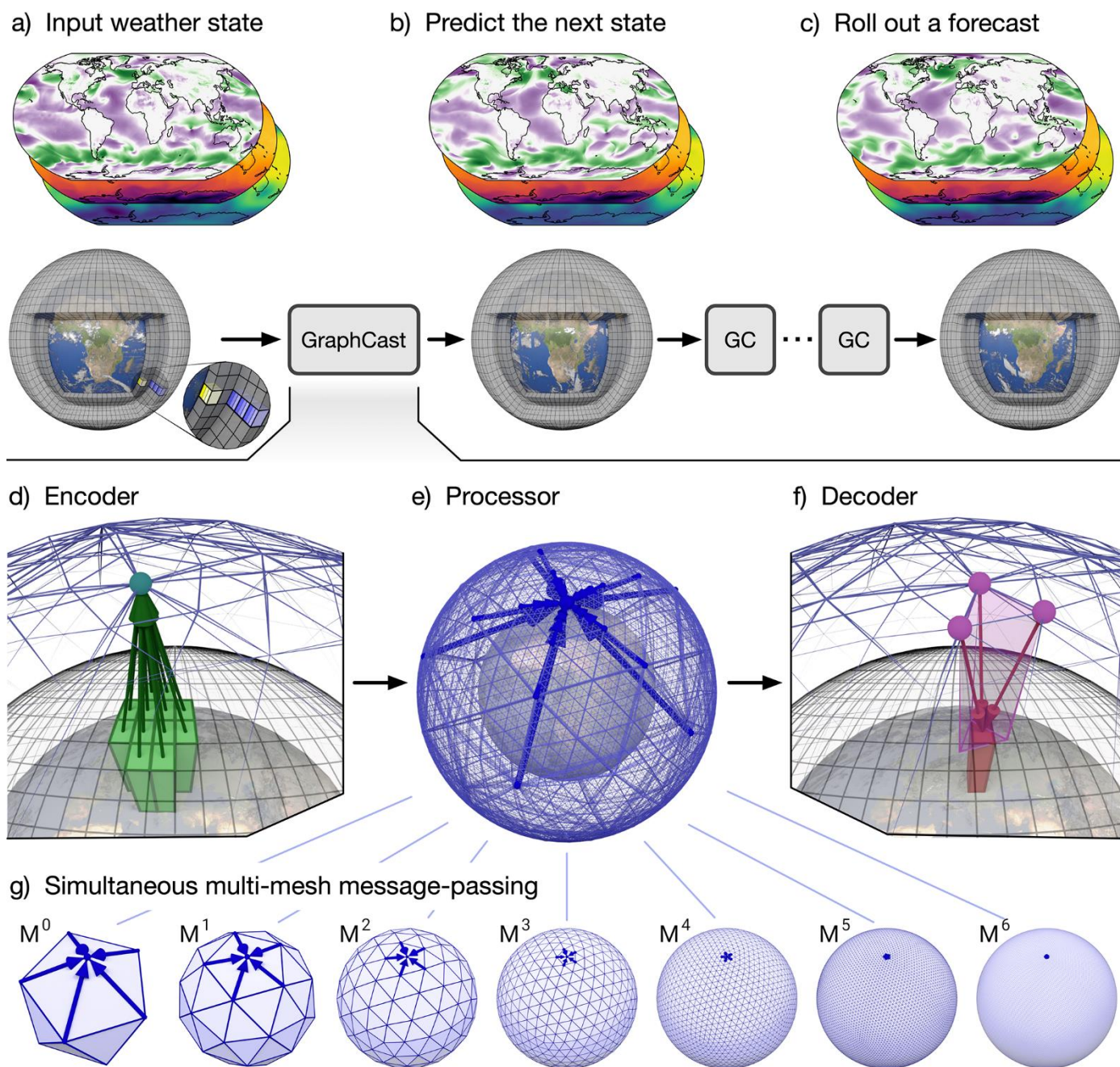
Published online 14 November 2023

10.1126/science.adi2336

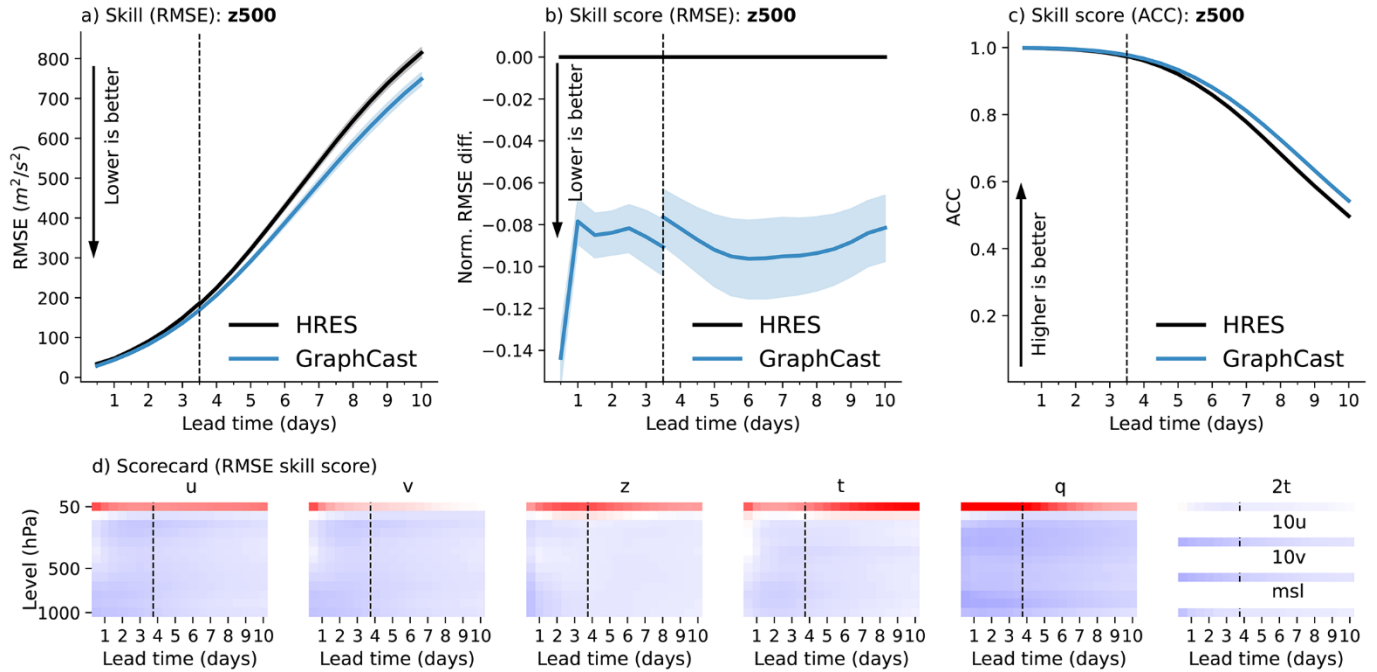
**Table 1. Weather variables and levels modeled by GraphCast.** The numbers in parentheses in the column headings are the number of entries in the column. Boldfaced variables and levels indicate those which were included in the scorecard evaluation. All atmospheric variables are represented at each of the pressure levels.

Surface variables (5)	Atmospheric variables (6)	Pressure levels (37)
<b>2-m temperature</b> (2 <i>T</i> )	<b>Temperature</b> ( <i>T</i> )	1, 2, 3, 5, 7, 10, 20, 30, <b>50</b> , 70,
<b>10 m u wind component</b> (10 <i>U</i> )	<b>U component of wind</b> ( <i>U</i> )	<b>100</b> , 125, <b>150</b> , 175, <b>200</b> , 225,
<b>10 m v wind component</b> (10 <i>V</i> )	<b>V component of wind</b> ( <i>V</i> )	<b>250</b> , <b>300</b> , 350, <b>400</b> , 450, <b>500</b> ,
<b>Mean sea-level pressure</b> ( <i>MSL</i> )	<b>Geopotential</b> ( <i>Z</i> )	550, <b>600</b> , 650, <b>700</b> , 750, 775,
Total precipitation ( <i>TP</i> )	<b>Specific humidity</b> ( <i>Q</i> )	800, 825, <b>850</b> , 875, 900, <b>925</b> ,
	Vertical wind speed ( <i>W</i> )	950, 975, <b>1000</b>



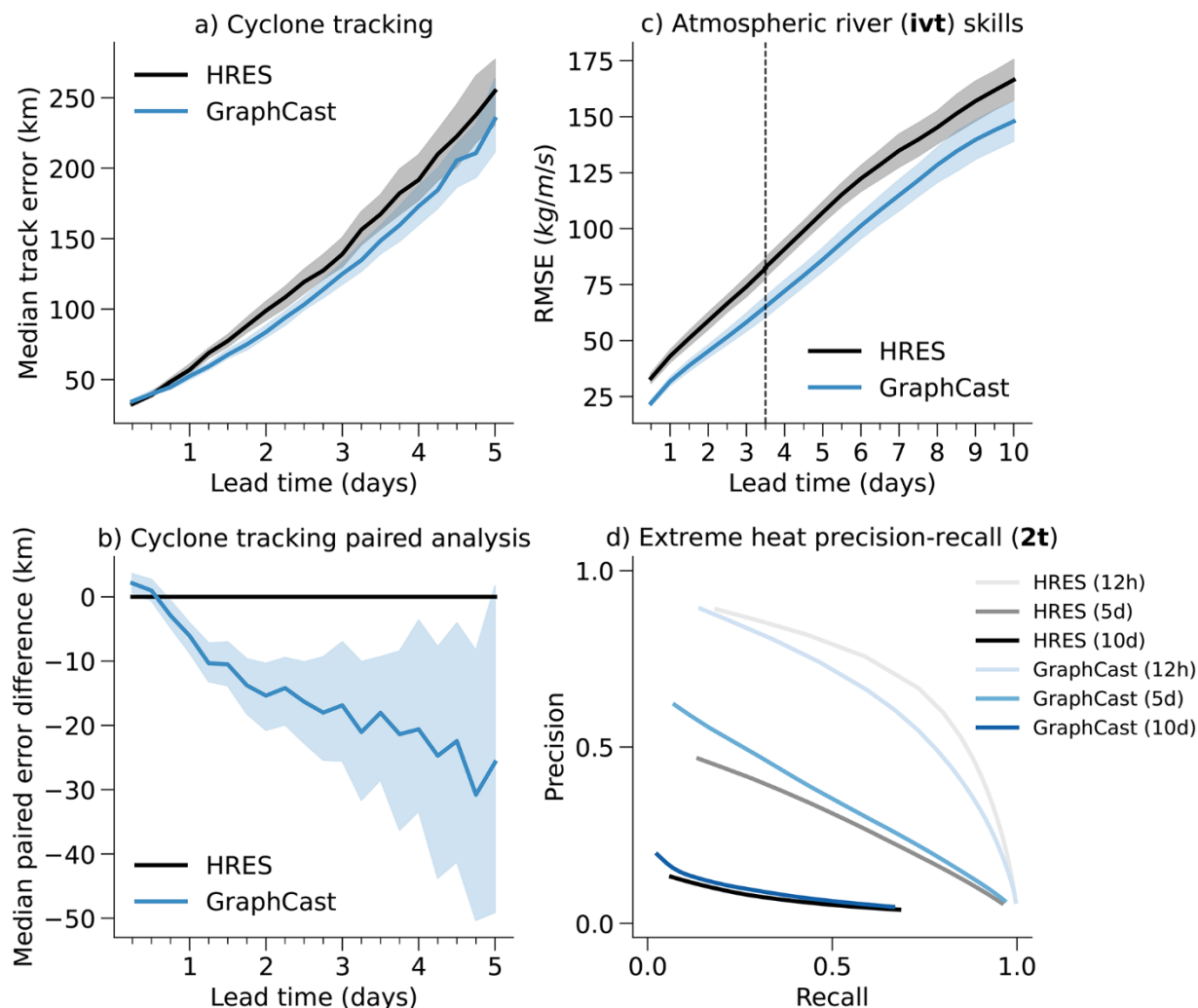


**Fig. 1. Model schematic.** (A) The input weather state(s) are defined on a  $0.25^\circ$  latitude-longitude grid comprising a total of  $721 \times 1440 = 1,038,240$  points. Yellow layers in the closeup pop-out window represent the 5 surface variables, and blue layers represent the 6 atmospheric variables that are repeated at 37 pressure levels ( $5 + 6 \times 37 = 227$  variables per point in total), resulting in a state representation of 235,680,480 values. (B) GraphCast predicts the next state of the weather on the grid. (C) A forecast is made by iteratively applying GraphCast to each previous predicted state, to produce a sequence of states which represent the weather at successive lead times. (D) The Encoder component of the GraphCast architecture maps local regions of the input (green boxes) into nodes of the multi-mesh graph representation (green, upward arrows which terminate in the green-blue node). (E) The Processor component updates each multi-mesh node using learned message-passing (heavy blue arrows that terminate at a node). (F) The Decoder component maps the processed multi-mesh features (purple nodes) back onto the grid representation (red, downward arrows which terminate at a red box). (G) The multi-mesh is derived from icosahedral meshes of increasing resolution, from the base mesh ( $M^0$ , 12 nodes) to the finest resolution ( $M^6$ , 40,962 nodes), which has uniform resolution across the globe. It contains the set of nodes from  $M^6$  and all the edges from  $M^0$  to  $M^6$ . The learned message-passing over the different meshes' edges happens simultaneously, so that each node is updated by all of its incoming edges. The Earth texture in the figure is used under CC~BY~4.0 from <https://www.solarsystemscope.com/textures/>.

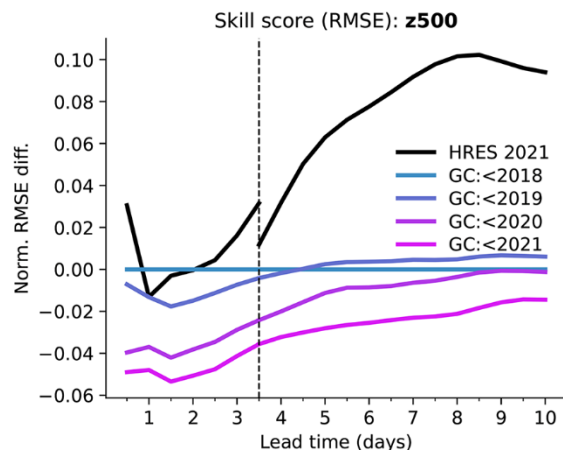


**Fig. 2. Global skill and skill scores for GraphCast and HRES in 2018.** (A) RMSE skill (y-axis) for GraphCast (blue lines) and HRES (black lines), on Z500, as a function of lead time (x-axis). Error bars represent 95% confidence intervals. The vertical dashed line represents 3.5 days, which is the last 12-hour increment of the HRES 06z/18z forecasts. The black line represents HRES, where lead times earlier and later than 3.5 days are from the 06z/18z and 00z/12z initializations, respectively. (B) RMSE skill score (y-axis) for GraphCast versus HRES, on Z500, as a function of lead time (x-axis). Error bars represent 95% confidence intervals for the skill score. We observe a discontinuity in GraphCast's curve because skill scores up to 3.5 days are computed between GraphCast (initialized at 06z/18z) and HRES's 06z/18z initialization, while after 3.5 days skill scores are computed with respect to HRES's 00z/12z initializations. (C) ACC skill (y-axis) for GraphCast (blue lines) and HRES (black lines), on Z500, as a function of lead time (x-axis). (D) Scorecard of RMSE skill scores for GraphCast, with respect to HRES. Each subplot corresponds to one variable:  $U$ ,  $V$ ,  $Z$ ,  $T$ ,  $Q$ ,  $2T$ ,  $10U$ ,  $10V$ ,  $MSL$ , respectively. The rows of each heatmap correspond to the 13 pressure levels (for the atmospheric variables), from 50 hPa at the top to 1000 hPa at the bottom. The columns of each heatmap correspond to the 20 lead times at 12-hour intervals, from 12 hours on the left to 10 days on the right. Each cell's color represents the skill score, as shown in (B), where blue represents negative values (GraphCast has better skill) and red represents positive values (HRES has better skill).





**Fig. 3. Severe-event prediction.** (A) Cyclone tracking performances for GraphCast and HRES. The x-axis represents lead times (in days), and the y-axis represents median track error (in km). Error bars represent bootstrapped 95% confidence intervals for the median. (B) Cyclone tracking paired error difference between GraphCast and HRES. The x-axis represents lead times (in days), and the y-axis represents median paired error difference (in km). Error bars represent bootstrapped 95% confidence intervals for the median difference (see supplementary materials section 8.1). (C) Atmospheric river prediction (IVT) skills for GraphCast and HRES. The x-axis represents lead times (in days), and the y-axis represents RMSE. Error bars are 95% confidence intervals. (D) Extreme heat prediction precision-recall for GraphCast and HRES. The x-axis represents recall, and the y-axis represents precision. The curves represent different precision-recall trade-offs when sweeping over gain applied to forecast signals (see supplementary materials section 8.3).



**Fig. 4. Training GraphCast on more recent data.** Each colored line represents GraphCast trained with data ending before a different year, from 2018 (blue) to 2021 (purple). The y-axis represents RMSE skill scores on 2021 test data, for Z500, with respect to GraphCast trained up to before 2018, over lead times (x-axis). The vertical dashed line represents 3.5 days, where the HRES 06z/18z forecasts end. The black line represents HRES, where lead times earlier and later than 3.5 days are from the 06z/18z and 00z/12z initializations, respectively.



## Learning skillful medium-range global weather forecasting

Remi Lam, Alvaro Sanchez-Gonzalez, Matthew Willson, Peter Wirsberger, Meire Fortunato, Ferran Alet, Suman Ravuri, Timo Ewalds, Zach Eaton-Rosen, Weihua Hu, Alexander Merose, Stephan Hoyer, George Holland, Oriol Vinyals, Jacklynn Stott, Alexander Pritzel, Shakir Mohamed, and Peter Battaglia

*Science* **Ahead of Print** DOI: 10.1126/science.adi2336

### View the article online

<https://www.science.org/doi/10.1126/science.adi2336>

### Permissions

<https://www.science.org/help/reprints-and-permissions>

Downloaded from <https://www.science.org> on December 20, 2023

Use of this article is subject to the [Terms of service](#)

---

*Science* (ISSN 1095-9203) is published by the American Association for the Advancement of Science. 1200 New York Avenue NW, Washington, DC 20005. The title *Science* is a registered trademark of AAAS.

Copyright © 2023 The Authors, some rights reserved; exclusive licensee American Association for the Advancement of Science. No claim to original U.S. Government Works