

EARLY ONLINE RELEASE

This is a PDF of a manuscript that has been peer-reviewed and accepted for publication. As the article has not yet been formatted, copy edited or proofread, the final published version may be different from the early online release.

This pre-publication manuscript may be downloaded, distributed and used under the provisions of the Creative Commons Attribution 4.0 International (CC BY 4.0) license. It may be cited using the DOI below.

The DOI for this manuscript is DOI:10.2151/jmsj.2025-031 J-STAGE Advance published date: June 26, 2025 The final manuscript after publication will replace the preliminary version at the above DOI once it is available.

| Intermediate Weather Forecasts among Multiple NWP |
|---|
| Models |
| |
| Atsushi KUDO ¹ |
| |
| Numerical Prediction Division |
| Japan Meteorological Agency, Tokyo, Japan |
| |
| |
| |
| |
| |
| |
| |
| November 15, 2024 |
| |
| |
| |
| |

Abstract

| 31 | Numerical weather prediction (NWP) centers around the world operate a variety of |
|----|--|
| 32 | NWP models. In addition, recent advances in AI-driven NWP models have further |
| 33 | increased the availability of NWP outputs. While this expansion holds the potential to |
| 34 | improve forecast accuracy, it raises a critical question: which prediction is the most |
| 35 | plausible? If the NWP models have comparable accuracy, it is impossible to determine in |
| 36 | advance which one is the best. Traditional approaches, such as ensemble or weighted |
| 37 | averaging, combine multiple NWP outputs to produce a single forecast with improved |
| 38 | accuracy. However, they often result in meteorologically unrealistic and uninterpretable |
| 39 | outputs, such as the splitting of tropical cyclone centers or frontal boundaries into multiple |
| 40 | distinct systems. |
| 41 | To address this issue, we propose DeepMedcast, a deep learning method that |
| 42 | generates intermediate forecasts between two or more NWP outputs. Unlike averaging, |
| 43 | DeepMedcast provides predictions in which meteorologically significant features—such |
| 44 | as the locations of tropical cyclones, extratropical cyclones, fronts, and shear lines- |
| 45 | approximately align with the arithmetic mean of the corresponding features predicted by |
| 46 | the input NWP models, without distorting meteorological structures. We demonstrate the |
| 47 | capability of DeepMedcast through case studies and verification results, showing that it |
| 48 | produces realistic and interpretable forecasts with higher accuracy than the input NWP |
| | 1 |

| 10 | modele. By providing plausible intermediate forecasts, Deepwededst our significantly |
|----|--|
| 50 | contribute to the efficiency and standardization of operational forecasting tasks, including |
| 51 | general, marine, and aviation forecasts. |

53 **Keywords** deep neural network; intermediate forecast; numerical weather prediction

55 **1. Introduction**

In recent decades, numerical weather predictions (NWPs) and their post-processing 56have played a central role in issuing weather forecasts, warnings, and advisories (WMO 572013; Vannitsem 2021). NWP centers around the world have developed and are operating 58a variety of NWP models for accurate weather predictions. For example, the European 59Centre for Medium-Range Weather Forecasts (ECMWF) operates the Integrated 60 Forecasting System (IFS) and its ensemble prediction system (ECMWF 2024); the UK Met 61Office operates the Unified Model and the Met Office Global and Regional Ensemble 62 Prediction System (Brown et al. 2012; Hagelin et al. 2017; Inverarity et al. 2023). The 63 National Centers for Environmental Prediction (NCEP) at the National Oceanic and 64 Atmospheric Administration (NOAA) operates the Global Forecast System (NCEP 2016), 65 the High-Resolution Rapid Refresh (Dowell et al. 2022), and the Hurricane Weather 66 Research and Forecasting model (Gopalakrishnan et al. 2011). The Japan Meteorological 67 Agency (JMA) operates three deterministic NWP models and two ensemble prediction 68 systems for short-range to weekly forecasts: the Global Spectrum Model (GSM), the Meso-69 70 Scale Model (MSM), the Local Forecast Model, the Global Ensemble Prediction System, and the Mesoscale Ensemble Prediction System (JMA 2024). These models cover different 71areas with varying resolutions and processes. 72

In addition to traditional physics-based NWP models, recent advancements in artificial
 intelligence (AI) have introduced new methods for producing weather predictions. Al-driven

For Peer Review

NWP models, such as FourCastNet (Pathak et al. 2022; Bonev et al. 2023), GraphCast
(Lam et al. 2022), Pangu-Weather (Bi et al. 2022; Bi et al. 2023), FengWu (Chen et al. 2023;
Han et al. 2024), Aurora (Bodnar et al. 2024), GenCast (Price et al. 2023), and AIFS (Lang
et al. 2024), have demonstrated the ability to enhance both the speed and accuracy of
weather predictions by leveraging deep learning techniques to model complex atmospheric
systems.

At present, forecasters are able to use multiple NWP models including AI-driven NWP 81 models, which provide a range of possible atmospheric states, allowing them to select the 82 most plausible prediction from available NWPs. However, this raises a critical question: 83 Which prediction is the most plausible? If the models have comparable accuracy, it is 84 impossible to determine in advance which one is the best. One practical and widely used 85solution is to average the results from multiple NWP models or their post-processed outputs, 86 as this can reduce random errors inherent in NWP models and can lead to higher accuracy 87 than individual models (Vislocky and Fritsch 1997; JMA 2018). The National Hurricane 88 Center and the Joint Typhoon Warning Center in the United States use consensus forecasts 89 (e.g., Simon et al. 2018; Cangialosi et al. 2023), which are weighted averages, extensively 90 for both tropical cyclone (TC) track and intensity predictions. JMA employs consensus 91forecasting for TC track predictions by averaging positions of TC centers from multiple NWP 92outputs to improve forecast accuracy (Nishimura and Fukuda 2019; JMA 2022). The UK Met 93 Office operates the IMPROVER system (Roberts et al. 2023), which applies a weighted 94

average of post-processing and nowcasts based on multiple NWP outputs. The National 95Weather Service (NWS) at NOAA operates the National Blend of Models (NBM), which 96 provides statistically post-processed multi-model ensemble guidance (Hamill et al. 2017). 97 The German Meteorological Service uses MOSMIX and ModelMIX (Primo et al. 2024), 98 which are weighted averages of post-processing based on IFS, their global model, and their 99 regional ensemble model. Additionally, the World Area Forecast Centre, comprising centers 100in London and Washington, operates harmonized forecasts, including mean, maximum, and 101 minimum forecasts, from both NWP outputs for aviation hazards such as cumulonimbus 102clouds, turbulence, and in-flight icing (ICAO 2016). 103104 It is straightforward to average the central position of TCs, extra tropical cyclones, or the location of fronts because averaging does not degrade their clarity. However, averaging 105atmospheric fields such as pressure or wind speed around these systems is not appropriate. 106This is because averaging can smooth out or distort these fields, weakening the central 107pressure or wind speeds around cyclones and fronts, and even causing TCs or fronts to split 108into two, resulting in predictions that are meteorologically unrealistic and difficult to interpret. 109110Forecasters must then choose between two options: using a single model that is realistic and interpretable but potentially less accurate or using an averaged prediction that is 111 unrealistic and uninterpretable but may be more accurate. 112Beyond these challenges, weather forecasting faces additional difficulties in operational 113

114 practice. In JMA's forecasting and warning issuance process, TC track forecasts, which are

| 115 | based on consensus forecasts derived by averaging the positions of TC centers from |
|-----|---|
| 116 | multiple NWP models, take precedence. As a result, forecasters responsible for general, |
| 117 | marine, and aviation forecasts must ensure that their forecasts align with TC track forecasts. |
| 118 | However, since no NWP model inherently conforms to the TC track forecasts, forecasters |
| 119 | need to adjust the existing NWP outputs in their minds to construct forecast scenarios that |
| 120 | align with them. This process requires significant time and effort and can pose a major |
| 121 | obstacle to the standardization of forecasting workflows, leading to inefficiencies in |
| 122 | operational forecasting. |
| 123 | In addition to these operational challenges, machine learning-based post-processing |
| 124 | presents its own set of difficulties, particularly regarding data requirements. In conventional |
| 125 | model output statistics (MOS), obtaining long-term, homogeneous datasets is particularly |
| 126 | difficult because the input NWP model is periodically updated, causing changes in its |
| 127 | systematic errors. Consequently, the statistical relationships learned from past data may no |
| | |

128 longer hold after a model update.

The objective of this study is to propose DeepMedcast, a method that uses deep learning to generate a realistic and interpretable "intermediate forecast" between two or more NWP models. In this study, we do not attempt to define intermediate forecasts in a physical or mathematical sense. Instead, we adopt a pragmatic definition: an intermediate forecast is a predicted meteorological field in which meteorologically significant features such as the center positions of TCs or extratropical cyclones, frontal boundaries, and shear

lines—are located at the arithmetic mean of the corresponding features predicted by input
 NWP models, and it simultaneously exhibits higher forecast accuracy against observations
 than the input models.

Unlike averaging, DeepMedcast can produce atmospheric fields around cyclones and 138fronts without smoothing out or disturbing their distributions. This capability is crucial in 139operational forecasts, where accurate and interpretable predictions are needed for issuing 140reliable warnings and advisories-particularly when a TC is approaching-and also 141contributes to the standardization of forecasting workflows by reducing reliance on manual 142adjustments by individual forecasters. While mathematical frameworks such as 143displacement interpolation and barycenters in optimal transport theory (McCann 1997; 144Agueh and Carlier 2011; Peyré and Cuturi, 2020) provide rigorously defined intermediate 145states and have been widely used in machine learning fields such as image processing, 146applying them directly to forecasts from multiple NWP models remains challenging. Le Coz 147et al. (2023) proposed a barycenter-based method using optimal transport to combine 148forecasts from multiple NWP models in the context of subseasonal prediction, offering a 149mathematically principled approach to multi-model ensemble forecasts. Duc and Sawada 150(2024) pointed out that traditional arithmetic ensemble means tend to excessively smooth 151rainfall structures, and proposed a Gaussian-Hellinger barycenter based on unbalanced 152optimal transport theory to derive more realistic and structurally coherent ensemble means. 153Their method is particularly effective in representing heavy precipitation and may contribute 154

For Peer Review

to future advances in ensemble post-processing and spatial verification techniques. To the 155best of the authors' knowledge, however, no prior study has successfully generated two-156dimensional intermediate forecasts for mean sea-level pressure or surface wind vectors at 157high temporal resolution and in areas strongly influenced by topography. In contrast, 158DeepMedcast is not intended to generate physically or mathematically consistent 159160intermediate forecasts but to provide forecasters with operationally practical solutions, addressing a critical challenge in operational forecasting. In many national meteorological 161centers, including JMA, machine learning-based post-processing methods, referred to as 162forecast guidance, are operationally employed. Although forecast guidance does not 163preserve physical consistency, it enhances forecast accuracy by reducing biases inherent 164in NWP models and significantly improves the efficiency of operational forecasting. 165DeepMedcast is a kind of post-processing method that is not designed to ensure physical 166consistency but to provide practical support for forecasters. 167

This paper is structured as follows: Section 2 presents the methodology and data used for DeepMedcast, detailing the deep learning architecture and training process. Section 3 discusses the results of applying DeepMedcast to multiple NWP models with case studies and verification results, and Section 4 offers a discussion of contributions to operational forecasting and the key features of DeepMedcast's architecture. Finally, Section 5 concludes with a summary of the findings and future research directions.

174

2. Method and data 175

2.1 The framework of DeepMedcast 176

The main idea behind DeepMedcast lies in its original approach to generating 177intermediate forecasts between two NWP outputs. Figure 1 illustrates the framework of 178Fig. 1 DeepMedcast. During the training phase, instead of using two different NWP outputs 179Fig. 2 intended for creating an intermediate forecast, DeepMedcast utilizes data at two forecast 180lead times (FT), FT = t - Δt and FT = t + Δt , from a single NWP model as input variables for 181 the deep neural network (DNN) (Fig. 1a). The output from the DNN is then compared to the 182forecast from the same NWP at the intermediate lead time (FT = t) to calculate the loss for 183the backpropagation process. This approach enables the network to generate intermediate 184forecasts while reducing the blurring effect often seen in machine learning (ML)-based post-185processing, as the input variables are not inherently affected by errors relative to observed 186values, and the predictions at FT = t are expected to lie between those at FT = t $\pm \Delta t$. During 187the inference phase, two different NWP outputs at the same forecast lead time are used to 188generate an intermediate forecast for the projection time (Fig. 1b). 189

190DeepMedcast is primarily designed to generate intermediate forecasts between two NWP models. However, the same DNN model can be applied recursively to generate 191intermediate forecasts between more than two NWP models. For instance, by taking 192intermediate forecasts between two pairs of NWP models, DeepMedcast can generate an 193intermediate result between four NWP models (Fig.2). This recursive approach could be 194

extended further to create intermediate forecasts between 8, 16, or even more NWP models.

196

197 2.2 Data used for the study

The NWP model used for training is GSM, which is operated by JMA four times a day 198(at 00, 06, 12, and 18 UTC as initial times). The training period spans nine years, from 199January 2013 to December 2021, while the validation period covers one year, from January 200to December 2022. GSM had a horizontal resolution of approximately 20 km until March 2012023, after which it was upgraded to about 13 km (JMA 2024). The GSM data used in this 202study is stored at JMA, where it is trimmed and linearly interpolated onto a 121 × 151 grid 203with a resolution of 0.25 degrees by 0.2 degrees around Japan (Fig. 3). Hereafter, we refer 204to this as the target grid domain. 205

The forecast variables include wind components (U, V), temperature (T), and relative 206humidity (RH) at both the surface and the 700 hPa level, as well as mean sea-level pressure 207at the surface (P_{sea}). To reduce computational cost and execution time, each variable is 208used individually to train separate networks, with each network dedicated to a single variable. 209210That means, each DNN model always takes two input channels (at FT = t - Δ t and FT = t + Δt) and outputs one channel (at FT = t) for DeepMedcast. Both input channels are utilized 211212by swapping their order, i.e., both FT = t - Δt and FT = t + Δt , as well as FT = t + Δt and FT = t - Δt , are employed to preserve symmetry. This strategy encourages the network to learn 213symmetric representations, so that meteorologically significant features—such as the center 214

positions of TCs or extratropical cyclones, frontal boundaries, and shear lines-215approximately align with the arithmetic mean of the corresponding features from the two 216input channels. The networks trained with 700 hPa data are employed to generate 217intermediate forecasts for the upper atmosphere (e.g., 850 hPa, 700 hPa, and 500 hPa) to 218reduce computational costs. The forecast lead times used in this study are t = 9, 10, 11, 12, 21913, and 14 hours, with $\Delta t = \pm 3$ and ± 6 hours, corresponding to t and Δt in Fig. 1, determined 220by taking both computational costs and accuracy into account. 221For the case studies in Section 3, MSM, IFS, GraphCast, and Pangu-Weather are used 222along with GSM for inference. MSM is operated by JMA eight times a day (at 00, 03, ..., and 22321 UTC as initial times) with a 5 km horizontal resolution, providing forecasts up to FT = 78 224hours for 00 and 12 UTC initial times and up to FT = 39 hours for other initial times. IFS 225data, provided by ECMWF for the World Meteorological Organization (WMO) members, has 226a horizontal resolution of 0.5 degrees and is initialized four times daily at 00, 06, 12, and 18 227UTC. Both GraphCast and Pangu-Weather have a horizontal resolution of 0.25 degrees, 228with data initialized at 00, 06, 12, and 18 UTC. These NWP outputs are linearly interpolated 229230to the target grid domain for inference.

231

232 2.3 DNN model architecture

In this study, a U-Net architecture (Ronneberger et al. 2015) is applied as the DNN
 model. The structure of the network is illustrated in Fig. 4. The encoder part of the U-Net

Fig. 4

For Peer Review

| 235 | utilizes a convolutional network with kernel size = 3, stride = 1, and padding = 1 for |
|-----|---|
| 236 | convolution operations. To progressively reduce the image size, MaxPooling layers with |
| 237 | kernel size = 2 and stride = 2 are employed. This downsampling process continues until the |
| 238 | image size is reduced to 1/8 of the original dimensions, at which point the channel count |
| 239 | reaches 2048, starting from an initial 2 channels that are expanded to 256 channels. In each |
| 240 | downsampling stage, the image size is halved while the number of channels doubles. In the |
| 241 | decoder part, transposed convolutional layers with kernel size = 2 and stride = 2 are applied |
| 242 | to upsample the feature maps, restoring the image to its original size while reducing the |
| 243 | channel count by half at each stage. By the final stage, the image is returned to its original |
| 244 | dimensions with 256 channels, which are then reduced to 1 channel in the output layer. The |
| 245 | activation functions used in this network include rectified linear units (ReLU, Nair and Hinton |
| 246 | 2010) for all layers except the output layer, which uses a sigmoid function to ensure output |
| 247 | values are scaled between 0 and 1. |

During the training phase, the two input channels and one ground truth channel, each consisting of 121×151 grids, are normalized to a value range of 0 to 1 using the maximum and minimum values across all three channels. Specifically, for T, RH, and P_{sea}, the normalization is applied as:

252

$$x' = \frac{x - x_{\min}}{x_{\max} - x_{\min}}$$

where x' is the normalized value, x is the input value, and x_{max} and x_{min} represent the maximum and minimum values, respectively. For wind components U and V, x_{max} and x_{\min} are defined as:

| 256 | $x_{\max} = \max(x_{\max} , x_{\min})$ |
|-----|---|
| 257 | $x_{\min} = -x_{\max}$ |
| 258 | and the same normalization is applied. |
| 259 | After normalization, the values are extended to 128 × 158 grids by copying the last |
| 260 | column and row to adjust to the network structure. The output values are compared with the |
| 261 | normalized and extended ground truth values using the mean square error (MSE) as the |
| 262 | loss function. We employ Adam (Kingma and Ba 2014) as optimization. |
| 263 | During the inference phase, the input values are normalized to the 0 to 1 range using |
| 264 | the maximum and minimum values of the two input channels. The output values are then |
| 265 | denormalized using the same maximum and minimum values, and resized back to 121 \times |
| 266 | 151 grids by trimming the extended columns and rows, providing predictions at the target |
| 267 | grid domain. |
| 268 | |
| 269 | 3. Results |
| 270 | In this section, we demonstrate the capability of DeepMedcast through four case studies |
| 271 | and verification results. The forecast data used here is from a period beginning in January |
| 272 | 2023, which is independent of the DNN's training and validation periods. The case studies |
| 273 | compare the atmospheric fields generated by DeepMedcast with those obtained via |
| 274 | arithmetic mean of the NWP outputs. |

3.1 Case 1: Position discrepancy in a typhoon forecast between GSM and MSM

The first case study examines a typhoon forecast where there is a positional discrepancy Fig. 5 between GSM and MSM. Figure 5 shows the predictions at FT = 51 hours based on the initial time of 12 UTC on 12 August 2023. This case focuses on Typhoon LAN which was moving northwest over the ocean south of Japan. At FT = 51 hours, GSM predicted the typhoon's position at 33.3°N, 137.1°E, while MSM placed it southwest at 32.7°N, 135.7°E. Both models predicted a central pressure of 960 hPa, with the maximum wind speed of 79 kt (1 kt \approx 0.514 m s⁻¹) (GSM) and 68 kt (MSM) (Figs. 5a and 5b).

When the mean sea-level pressure and surface wind components from GSM and MSM were averaged arithmetically, the typhoon's center split into two, aligning with the predicted positions from each model (Fig. 5c). Such a result is evidently unrealistic and lacks interpretability. The central pressure weakened to 974 hPa, and the maximum wind speed reduced to 58 kt, which made the forecast meteorologically unnatural, with the typhoon taking on an elongated structure and weakening wind speeds near the center. This resulted in a forecast that was difficult to explain and potentially misleading.

In contrast, DeepMedcast generated a plausible forecast, placing the typhoon at 33.0°N, 136.4°E, halfway between GSM and MSM predictions (Fig. 5d). The typhoon maintained a single, natural, and interpretable shape with a central pressure of 960 hPa and the maximum wind speed of 69 kt, representing an intermediate forecast between the two NWP models.

Fig. 6

Fig. 7

Figure 6 compares the DeepMedcast outputs with different input orders. Figure 6a 295shows the result when GSM and MSM are used as inputs in that order (identical to Fig. 5d), 296while Fig. 6b shows the result when the input order is reversed (MSM-GSM). As expected, 297the predicted structure-including the typhoon center position, central pressure, and 298surrounding wind field-remains qualitatively consistent, despite minor differences due to 299the asymmetry of the trained neural network. This indicates that while DeepMedcast is not 300 strictly order-invariant, the resulting intermediate forecasts are robust to changes in input 301 order. 302

303

304 3.2 Case 2: Discrepancy in a front position forecast between GSM and MSM

The second case study examines a forecast where there was a positional discrepancy 305 in the predicted location of a front between GSM and MSM. Figure 7 shows the predictions 306 at FT = 30 hours based on the initial time of 00 UTC on 17 June 2024. At the initial time, a 307 stationary front was located south of Japan (not shown), and by FT = 30 hours, the front 308 was predicted to move northward toward Tokyo (indicated by the blue circles in the figure). 309 GSM predicted the front to the south of Tokyo, indicated by the blue dashed line, with a 310clear wind direction and speed shear, which corresponds well with the 21°C isotherm around 311Tokyo (Fig. 7a). In contrast, MSM placed the front north of Tokyo (green dashed line), also 312with a clear wind direction and speed shear aligned with the 21°C isotherm (Fig. 7b), 313resulting in a positional discrepancy between the two NWP models. Consequently, GSM 314

Page 17 of 100

predicted a northerly to northeasterly wind and cooler temperatures around Tokyo, while
 MSM predicted southerly to southwesterly winds and warmer temperatures, leading to
 significant differences in the forecast for Tokyo.

The arithmetic mean of the GSM and MSM predictions (Fig. 7c) results in a split structure for the front, with wind shear corresponding to the locations predicted by GSM and MSM (shown by the purple lines), while the 21°C isotherm is predicted between the two fronts. This demonstrates that when there is a discrepancy in the predicted front position, simple averaging of the atmospheric fields leads to an unnatural and uninterpretable forecast that cannot maintain the original front structure.

In contrast, DeepMedcast (Fig. 7d) generates a clear wind direction and speed shear at the intermediate position between the GSM and MSM predictions (indicated by the brown dashed line), which corresponds well with the 21°C isotherm. DeepMedcast successfully produces a realistic and interpretable intermediate forecast while preserving the structure of the original front.

329

330 3.3 Case 3: Significant difference in low-pressure system position between GSM and MSM

Fig. 8

The third case study highlights a situation where there was a large difference in the predicted position of a low-pressure system between GSM and MSM. Figure 8 shows surface wind and mean sea-level pressure, along with temperature and dew-point depression at 850 hPa, at FT = 75 hours based on the initial time of 12 UTC on 28 August

2024. At the initial time, Typhoon SHANSHAN was located south of Kyushu (see Fig. 3, location 4) at 30.6°N, 130.2°E, slowly moving northward (not shown). By 15 UTC on 31 August (FT = 75 hours), the system, which had either remained a tropical storm or transitioned into a low-pressure system, was predicted by GSM to be south of the Kanto region (see Fig. 3, location 3) at 35.4°N, 139.8°E (Fig. 8a), while MSM placed it east of Hokkaido (see Fig. 3, location 1) at 42.4°N, 147.3°E (Fig. 8b).

This case highlights a large positional difference of about 1000 km between the GSM 341and MSM predictions. When the arithmetic mean of these is taken (Fig. 8c), it results in two 342distinct low-pressure systems at the positions predicted by each model, creating an 343uninterpretable forecast. In contrast, DeepMedcast predicted a single low-pressure system 344located between the two forecasts, at 39.5°N, 143.0°E, off the Pacific coast of Tohoku (see 345Fig. 3, location 2; Fig. 8d). Additionally, when examining the moisture area at 850 hPa (dew-346point depression < 3°C), the arithmetic mean shows moist areas surrounding both the GSM 347and MSM low-pressure systems, with a relatively dry region in between. On the other hand, 348DeepMedcast represents a moist area around its low-pressure system, corresponding well 349 with the surface pressure field, providing a realistic and interpretable forecast. 350

351

352 3.4 Case 4: Intermediate forecast between four NWP models for typhoon KHANUN

353 The fourth case study presents an intermediate forecast between four NWP models:

354 GSM, IFS, GraphCast, and Pangu-Weather. Figure 9 shows surface wind and mean sea-

Fig. 9

Table.1

For Peer Review

| 355 | level pressure at FT = 108 hours, based on the initial time of 12 UTC on 2 August 2023. At |
|-----|---|
| 356 | the initial time, Typhoon KHANUN was located west of Okinawa (see Fig. 3, location 5) at |
| 357 | 26.2°N, 125.6°E, slowly moving westward (not shown). By 00 UTC on 7 August (FT = 108 |
| 358 | hours), the central position was predicted by the four models to be at 31.4°N, 131.1°E (GSM), |
| 359 | 28.9°N, 133.0°E (IFS), 28.0°N, 131.2°E (GraphCast), and 28.8°N, 130.7°E (Pangu- |
| 360 | Weather). The typhoon central position, central pressure, and maximum wind speed at FT |
| 361 | = 108 hours for each model are summarized in Table 1. |

When the arithmetic mean of the four models is taken, the center splits into two, with a 362 weakened central pressure of 979 hPa and the maximum wind speed of 36 kt (Fig. 9e), both 363the same or weaker than the predictions of any individual model. In contrast, DeepMedcast 364365 predicted a single center at 29.3°N, 131.4°E, with the central pressure of 964 hPa and the maximum wind speed of 42 kt (Fig. 9f), representing an intermediate intensity forecast 366 between the four NWP models. The average central latitude, longitude, and pressure of the 367 four NWP models were 29.3°N, 131.5°E, and 964 hPa, respectively, matching 368 DeepMedcast's prediction. The average maximum wind speed of the four models was 49 369 370 kt, meaning DeepMedcast's forecast was slightly weaker than the average.

Figure 10 illustrates the effect of changing the order in which intermediate forecasts are taken when combining the four NWP models. Figure 10a is identical to Fig. 9f and shows the result when intermediate forecasts are first generated between GSM and IFS, and between GraphCast and Pangu-Weather, followed by taking an intermediate forecast

Fig. 10

375between those two results. Figure 10b shows the case where the intermediate forecasts are first taken between GSM and GraphCast, and between IFS and Pangu-Weather, then 376 combined. Figure 10c presents the result when intermediate forecasts are first taken 377 between GSM and Pangu-Weather, and between IFS and GraphCast. As in the two-model 378case discussed in Section 3.1, the outputs differ slightly due to the inherent asymmetry of 379 the trained network and the recursive nature of the procedure. However, the predicted 380typhoon structure, including its central position, pressure, and wind field, remains 381 qualitatively consistent across all three cases. This suggests that although DeepMedcast is 382not strictly order-invariant, it produces robust intermediate forecasts in practice. 383

384

385 3.5 Statistical evaluation of DeepMedcast using surface wind observations

This section presents the verification results of DeepMedcast generated from GSM and MSM. Surface wind predictions from DeepMedcast, along with the input GSM and MSM forecasts, were verified against observations from the Automated Meteorological Data Acquisition System (AMeDAS), an automated observation network operated by JMA. The verification metric is the root mean square error (RMSE), defined as follows:

391
$$RMSE = \sqrt{\frac{1}{T} \sum_{t=1}^{T} \frac{1}{N} \sum_{n=1}^{N} (F_{nt} - O_{nt})^2}$$

where *T* and *N* are the numbers of forecast times and stations used for verification, and F_{nt} and O_{nt} represent the forecast and observed winds at station *n* and time *t*,

| 394 | respectively. Predictions from DeepMedcast, GSM, and MSM were linearly interpolated to |
|-----|--|
| 395 | each AMeDAS station from the four surrounding grid points. |

| 396 | Figure 11 shows the RMSE of wind speed (Fig. 11a) and wind direction (Fig. 11b) by | Fig. 11 |
|-----|--|---------|
| 397 | forecast lead time, ranging from 3 to 39 hours. The verification was conducted over one | |
| 398 | year, from January to December 2023, using forecasts initialized four times daily (00, 06, | |
| 399 | 12, and 18 UTC), independent of the training and validation periods. In both panels, the red, | |
| 400 | blue, and green lines represent DeepMedcast, GSM, and MSM, respectively. As shown in | |
| 401 | Fig. 11, DeepMedcast achieves lower RMSEs for both wind speed and direction across all | |
| 402 | forecast lead times compared to its input models. The RMSE for wind direction in Fig. 11b | |
| 403 | shows a 6-hourly fluctuation pattern, reflecting the four-times-daily initialization and diurnal | |
| 404 | variation. | |

406 **4. Discussion**

407 4.1 Contributions to operational forecasting

As demonstrated by the case studies and verification in Section 3, DeepMedcast is capable of generating plausible and interpretable intermediate forecasts. The ability to generate intermediate forecasts between multiple models is expected to significantly contribute to operational forecasting. As mentioned in the Introduction, TC track forecasts, which are based on consensus from multiple NWP models, specifically the average position of the TC center predicted by these models, serve as primary reference in JMA's operational forecasting. Consequently, forecasters responsible for general, marine, and aviation forecasts must ensure that their forecasts align with the TC track forecasts. However, since no NWP model inherently conforms to the TC track forecasts, forecasters must adjust the existing NWP outputs in their minds to construct forecast scenarios that follow the TC track forecasts. DeepMedcast has a capability to provide two-dimensional wind and pressure fields that align with TC track forecasts, which could greatly improve the efficiency and standardization of tasks for operational forecasting.

Additionally, DeepMedcast can be effectively utilized in operational forecasting when 421there are discrepancies in the predicted positions of low-pressure systems or fronts among 422multiple models. When significant differences exist among NWP models, forecasters need 423to choose between two options: either using one model as the main scenario while treating 424others as alternative scenarios, or applying averaging methods. By generating an 425intermediate state between two or more NWP models, DeepMedcast provides a forecast 426 scenario that is more plausible than individual NWP models. However, it should be noted 427that it does not explicitly represent the variability among the original NWP models. This issue 428429is not unique to DeepMedcast but is also present when using arithmetic or weighted averaging of multiple NWP models or post-processing techniques. Since variability is an 430 indicator of forecast uncertainty, especially for longer lead times, it is important to take this 431information into account in daily operational forecasting. One practical approach is to 432incorporate the values of the original NWP models and their spread alongside the 433

intermediate forecast. This framework can further enhance operational forecasting by
 offering a practical way to take model variability into account.

436

437 4.2 Key Features of DeepMedcast's Architecture

There are two important features associated with DeepMedcast's architecture. The first 438is its flexibility in increasing the amount of training data. As mentioned in Section 2, this study 439used t = 9–14 and Δt = ±3, ±6, and in our experience, increasing t and Δt leads to better 440 forecast representation. One common issue in training DNN models is a lack of sufficient 441training data (e.g., Deng 2009; LeCun 2015). However, in the case of DeepMedcast, more 442training data can easily be generated by increasing t and Δt or by adding additional NWP 443444models. It is important to note, though, that increasing t and Δt requires more memory and computational time, which should be considered when expanding the dataset. 445

The second is DeepMedcast's maintainability. Despite being trained solely on 20-km 446 resolution GSM data, DeepMedcast works effectively not only with 13-km resolution GSM 447data but also with other NWP models such as MSM, IFS, GraphCast, and Pangu-Weather. 448 449 This is significant because most AI or ML methods in meteorology learn the relationship between input and target data, and when the characteristics of the input data change due to 450NWP model updates, retraining, fine-tuning, and/or online learning are usually required. 451While this is an unavoidable task for most AI or ML methods in meteorology, it is a time-452consuming yet essential task that operational centers have traditionally managed. However, 453

454 DeepMedcast can be applied to various NWP models without updating the DNN model since
455 it is not designed to correct NWP model biases, which significantly reduces the maintenance
456 costs for operational centers.

457

458 **5.** Summary

In this paper, we introduced DeepMedcast, a novel deep learning-based approach for producing intermediate forecasts between two or more NWP models. DeepMedcast was developed to generate plausible and interpretable intermediate forecast, bridging the gap between NWP model outputs.

A key advantage of DeepMedcast is its applicability to various NWP outputs without the need for retraining or fine-tuning the DNN. By providing plausible intermediate forecasts, DeepMedcast can significantly enhance the efficiency and standardization of operational forecasting tasks, including general, marine, and aviation forecasts.

Although DeepMedcast introduces some advancements, further research and development are needed to address several challenges. In this study, U-Net was employed as the DNN architecture; however, advanced methods such as Transformers (Vaswani et al. 2017; Dosovitskiy et al. 2020) and Diffusion models (Song and Ermon 2019; Ho et al. 2020) may further enhance DeepMedcast's representational capabilities. As shown in the case study in Section 3.4, the current method tends to slightly underestimate the maximum wind speed near TCs. Enhancing the DNN could help resolve this issue. Additionally, while

For Peer Review

474this study trained separate networks for each physical variable to reduce computational cost, incorporating multiple physical variables as input could potentially enhance forecast 475accuracy. By developing post-processing methods that use DeepMedcast as input, it would 476be possible to provide even more accurate predictions. This study focused on generating 477intermediate forecasts using a 1:1 weighting ratio, meaning that the two input models were 478given equal weight. Future work should explore methods for generating intermediate 479forecasts with other weighting ratios, such as 1:2. This would enable the application of 480 DeepMedcast to cases involving several models that is not a power of two, such as finding 481an intermediate forecast among three models. Furthermore, while this study demonstrated 482intermediate forecasts using two or four NWP models, DeepMedcast could be extended to 4838, 16, or more inputs, enabling the use of multiple NWP and ensemble models. Lastly, this 484study did not address intermediate precipitation forecasts. Since precipitation is one of the 485most critical variables in weather forecasting, future work will focus on developing 486intermediate precipitation forecasts. 487

488

489 Data Availability Statement

The datasets used in this study are available from the following sources: The GSM and MSM 490 data are operationally produced by JMA and can be accessed through the Japan 491Meteorological Business Support Center (http://www.jmbsc.or.jp/en/index-e.html). The IFS 492WMO ECMWF website 493data are accessible to members via the

(https://www.ecmwf.int/en/forecasts/datasets/wmo-additional). The GraphCast and PanguWeather source code and plugins are available under open-source licenses in the ECMWF
GitHub repository (https://github.com/ecmwf-lab/ai-models). Pre-trained models of PanguWeather and GraphCast, used without modification to generate forecast data, are
specifically accessible at https://github.com/ecmwf-lab/ai-models-panguweather and
https://github.com/ecmwf-lab/ai-models-graphcast.

- 500
- 501

Acknowledgments

We acknowledge ECMWF for providing IFS data available to WMO members through their website (https://www.ecmwf.int/en/forecasts/datasets/wmo-additional). Additionally, we are grateful for the availability of GraphCast and Pangu-Weather and extend our thanks to their respective developers. GraphCast and Pangu-Weather were used without modification to generate forecast data and are accessible via the ECMWF AI GitHub repository (https://github.com/ecmwf-lab/ai-models), supported by ECMWF. The author declares no conflicts of interest associated with this manuscript.

- 509
- 510

References

Agueh, M. and G. Carlier, 2011: Barycenters in the Wasserstein space SIAM Journal on

512 Mathematical Analysis, 43(2), 904–924. [Available at https://hal.science/hal-

513 00637399v1/document]

523

| 514 | Bi, K., L. Xie, H. Zhang, X. Chen, X. Gu, and Q. Tian, 2022: Pangu-Weather: A 3D High- |
|-----|---|
| 515 | Resolution System for Fast and Accurate Global Weather Forecast. arXiv preprint |
| 516 | arXiv:2211.02556. |
| 517 | Bi, K., L. Xie, H. Zhang, X. Chen, X. Gu, and Q. Tina, 2023: Accurate medium-range |
| 518 | global weather forecasting with 3D neural networks. <i>Nature</i> , 619 , 533–538. |
| 519 | https://doi.org/10.1038/s41586-023-06185-3. |
| 520 | Bodnar, C., W. P. Bruinsma, A. Lucic, M. Stanley, J. Brandstetter, P. Garvan, M. Riechert |
| 521 | J. Weyn, H. Dong, A. Vaughan, J. K. Gupta, K. Tambiratnam, A. Archibald, E. Heider, |
| | |

Atmosphere. arXiv preprint arXiv:2405.13063.

Bonev, B., T. Kurth, C. Hundt, J. Pathak, M. Baust, K. Kashinath, A. Anandkumar, 2023: 524

M. Welling, R. E. Turner, P. Perdikaris. 2024: Aurora: A Foundation Model of the

Spherical Fourier Neural Operators: Learning Stable Dynamics on the Sphere. arXiv 525

preprint arXiv:2306.03838. 526

Brown, A., S. Milton, M. Cullen, B. M. J. Golding, and A. Shelly, 2012: Unified modeling 527

and prediction of weather and climate: a 25 year journey, Bull. Amer. Meteor. Soc., 93, 528

5291865-1877.

Cangialosi, J., B.J. Reinhart, and J. Martinez, 2023: National Hurricane Center verification 530

report, 2023 Hurricane Season. Natinal Hurricane Center, 81pp. [Available at 531

https://www.nhc.noaa.gov/verification/pdfs/Verification_2023.pdf.] 532

Chen, K., T. Han, J. Gong, L. Bai, F. Ling, J. Luo, X. Chen, L. Ma, T. Zhang, R. Su, Y. Ci, 533

- B. Li, X. Yang, W. Ouyang, 2023: FengWu: Pushing the Skillful Global Medium-range
- 535 Weather Forecast beyond 10 Days Lead. *arXiv preprint arXiv:2304.02948.*
- 536 Deng, J., W. Dong, R. Socher, L. Li, K. Li, and L. Fei-Fei, 2009: ImageNet: A large-scale
- 537 hierarchical image database. IEEE Conference on Computer Vision and Pattern
- 538 Recognition, 248–255.
- 539 Dosovitskiy, A., L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M.
- 540 Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby, 2020: An
- 541 Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. *arXiv*
- 542 preprint arXiv:2010.11929.
- 543 Dowell, D. C., C. R. Alexander, E. P. James, S. S. Weygandt, S. G. Benjamin, G. S.
- 544 Manikin, B. T. Blake, J. M. Brown, J. B. Olson, M. Hu, T. G. Smirnova, T. Ladwig, J. S.
- 545 Kenyon, R. Ahmadov, D. D. Turner, J. D. Duda, T. I. Alcott, 2022: The High-Resolution
- 546 Rapid Refresh (HRRR): An hourly updating convection-allowing forecast model. Part I:
- 547 Motivation and system description. *Wea. Forecasting*, **37**, 1371–1395,
- 548 https://doi.org/10.1175/WAF-D-21-0151.1.
- 549 Duc, L. and Y. Sawada, 2024: Geometry of rainfall ensemble means: from arithmetic
- averages to Gaussian-Hellinger barycenters in unbalanced optimal transport. J. Meteor.
- 551 Soc. Japan, **102**, 35-47, https://doi.org/10.2151/jmsj.2024-003.
- 552 ECMWF, 2024: IFS Documentation. European Centre for Medium-Range Weather
- 553 Forecasts. [Available at https://www.ecmwf.int/en/publications/ifs-documentation.]

| 554 | Gopalakrishnan, S. G., F. Marks Jr., X. Zhang, JW. Bao, KS. Yeh, and R. Atlas, 2011: |
|-----|--|
| 555 | The experimental HWRF system: A study on the influence of horizontal resolution on the |
| 556 | structure and intensity changes in tropical cyclones using an idealized framework. Mon. |
| 557 | Wea. Rev., 139 , 1762–1784. |
| 558 | Hagelin, S., J. Son, R. Swinbank, A. McCabe, N. Roberts, and W. Tennant, 2017: The Met |
| 559 | Office convective-scale ensemble, MOGREPS-UK. Quart. J. Roy. Meteor. Soc., 143, |
| 560 | 2846–2861. |
| 561 | Hamill, T. M., E. Engle, D. Myrick, M. Peroutka, C. Finan, and M. Scheuerer, 2017: The |
| 562 | U.S. National Blend of Models for Statistical Postprocessing of Probability of |
| 563 | Precipitation and Deterministic Precipitation Amount. Mon. Wea. Rev., 145, 3441–3463. |
| 564 | Han T., S. Guo, F. Ling, K. Chen, J. Gong, J. Luo, J. Gu, K. Dai, W. Ouyang, L. Bai, 2024: |
| 565 | FengWu-GHR: Learning the Kilometer-scale Medium-range Global Weather |
| 566 | Forecasting. <i>arXiv preprint</i> arXiv:2402.00059. |
| 567 | Ho, J., A. Jain, and P. Abbeel, 2020: Denoising Diffusion Probabilistic Models. arXiv |
| 568 | preprint arXiv:2006.11239. |
| 569 | ICAO, 2016: Guidance on the harmonized WAFS grids for Cumulonimbus cloud, icing and |
| 570 | turbulence forecasts (version 2.6). International Civil Aviation Organization, 16pp. |
| 571 | [Available at |
| 572 | https://www.icao.int/airnavigation/METP/MOG%20WAFS%20Reference%20Documents/ |
| 573 | WAFS_HazardGridUserGuide.pdf.] |
| | 28 |

| 574 | Inverarity, G | G. W., | W. J. ⁻ | Tennant, I | L. Anton, | N. E. I | Bowler, | A. M. | Clayton, | M. Jaro | dak, A | . С |
|-----|---------------|--------|--------------------|------------|-----------|---------|---------|-------|----------|---------|--------|-----|
|-----|---------------|--------|--------------------|------------|-----------|---------|---------|-------|----------|---------|--------|-----|

- Lorence, F. Rawlins, S. A. Thompson, M. S. Thurlow, D. N. Walters, and M. A. Wlasak,
- 576 2023: Met Office MOGREPS-G initialisation using an ensemble of hybrid four-
- dimensional ensemble variational (En-4DEnVar) data assimilations. Quart. J. Roy.
- 578 *Meteor. Soc.*, **149**, 1138–1164.
- 579 JMA, 2018: Instruction for guidance. *Report of Numerical Prediction Division*. 64, Japan
- 580 Meteorological Agency, 248 pp (in Japanese). [Available at
- 581 https://www.jma.go.jp/jma/kishou/books/nwpreport/64/No64_all.pdf.]
- JMA, 2022: Annual Report on the Activities of the RSMC Tokyo Typhoon Center 2022.
- Japan Meteorological Agency, 143pp. [Available at https://www.jma.go.jp/jma/jma-
- eng/jma-center/rsmc-hp-pub-eg/AnnualReport/2022/Text/Text2022.pdf.]
- 585 JMA, 2024: Outline of the operational numerical weather prediction at the Japan
- 586 *Meteorological Agency.* Japan Meteorological Agency, 262pp. [Available at
- 587 https://www.jma.go.jp/jma/jma-eng/jma-center/nwp/outline-latest-nwp/index.htm.]
- 588 Kingma, D. P., and J. Ba, 2014: Adam: A Method for Stochastic Optimization. arXiv
- 589 *preprint arXiv:1412.6980.*
- Lam, R., A. Sanchez-Gonzalez, M. Willson, P. Wirnsberger, M. Fortunato, F. Alet, S.
- Ravuri, T. Ewalds, Z. Eaton-Rosen, W. Hu, A. Merose, S. Hoyer, G. Holland, O. Vinyals,
- J. Stott, A. Pritzel, S. Mohamed and P. Battaglia, 2022: GraphCast: Learning skillful
- 593 medium-range global weather forecasting. *arXiv preprint arXiv:2212.12794*.

- Lang, S., M. Alexe, M. Chantry, J. Dramsch, F. Pinault, B. Raoult, M. C. A. Clare, C.
- Lessig, M. Maier-Gerber, L. Magnusson, Z. B. Bouallègue, A. P. Nemesio, P. D.
- 596 Dueben, A. Brown, F. Pappenberger, F. Rabier, 2024: AIFS ECMWF's data-driven
- forecasting system. *arXiv preprint arXiv:2406.01465*.
- 598 LeCun, Y., Y. Bengio, and G. Hinton, 2015: Deep learning. *Nature*, **521**, 436–444.
- Le Coz, C., A. Tantet, R. Flamary, and R. Plougonven, 2023: A barycenter-based
- approach for the multi-model ensembling of subseasonal forecasts. arXiv preprint
- 601 *arXiv:2310:17933*.
- McCann, R. J., 1997: A Convexity Principle for Interacting Gases. Advances in
- 603 *Mathematics*, **128(1)**, 153–179.
- Nair, V., and G. E. Hinton, 2010: Rectified linear units improve restricted Boltzmann
- machines. Proceedings of the Twenty-seventh International Conference on Machine
- Learning (ICML-10), Haifa, Israel, 807–814.
- NCEP, 2016: Global Forecast System Global Spectral Model (GSM) v13.0.2. [Available
- at https://vlab.noaa.gov/web/gfs/documentation.]
- Nishimura, S., and J. Fukuda, 2019: Advancement of Tropical Cyclone Track Forecasts.
- 610 Yohou Gijutsu Kenshu Text, 24, 114–141 (in Japanese). [Available at
- 611 https://www.jma.go.jp/jma/kishou/books/yohkens/24/all.pdf.]
- Pathak, J., S. Subramanian, P. Harrington, S. Raja, A. Chattopadhyay, M. Mardani, T.
- Kurth, D. Hall, Z. Li, K. Azizzadenesheli, P. Hassanzadeh, K. Kashinath, and A.

- 614 Anandkumar, 2022: FourCastNet: A Global Data-driven High-resolution Weather Model
- using Adaptive Fourier Neural Operators. *arXiv preprint arXiv:2202.11214*.
- 616 Peyré, G., and M. Cuturi, 2020: Computational Optimal Transport. arXiv preprint arXiv:

617 **1803.00567**

- ⁶¹⁸ Price, I., A. Sanchez-Gonzalez, F. Alet, T. R. Andersson, A. El-Kadi, D. Masters, T.
- Ewalds, J. Stott, S. Mohamed, P. Battaglia, R. Lam, M. Willson, 2023: GenCast:
- Diffusion-based ensemble forecasting for medium-range weather. *arXiv preprint*
- 621 arXiv:2312.15796.
- Primo, C., B. Schulz, S. Lerch, and R. Hess, 2024: Comparison of Model Output Statistics
- and Neural Networks to Postprocess Wind Gusts. *arXiv preprint arXiv:2401.11896*.
- Roberts, N., A. Benjamin, E. Gavin, M. Stephen, R. Fiona, S. Caroline, T. Tomasz, A.
- Paul, B, Laurence. C. Neil, F. Ben, F. Jonathan, G. Tom, H. Leigh, H. Aaron, H.
- Katharine, J. Simon, J. Caroline, M. Ken, S. Christopher, S. Michael, W. Bruce, B.
- Simon, B. Mark, B. Daniel, B. Anna, B. Clare, C. Robert, C. Sean, C. Ric, H. Roger, H.
- Kathryn, H. Teresa, M. Marion, P. Jon, P. Tim, S. Victria, S. Eleanor, and W. Mark,
- 2023: IMPROVER: The New Probabilistic Postprocessing System at the Met Office. Bull
- 630 Amer. Meteor. Soc., **104**, E680–E697.
- Ronneberger, O., P. Fischer, and T. Brox, 2015: U-Net: Convolutional Networks for
- Biomedical Image Segmentation. *arXiv preprint arXiv:1505.04597*.
- 633 Simon, A., A. B. Penny, M. DeMaria, J. L. Franklin, R. J. Pasch, E. N. Rappaport, and D.

- A. Zelinsky, 2018: A description of the real-time HFIP Corrected Consensus Approach
- 635 (HCCA) for Tropical Cyclone Track and Intensity Guidance. Wea. Forecasting, 33, 37–
- 636 **57**, https://doi.org/10.1175/WAF-D-17-0068.1.
- 637 Song, Y. and S. Ermon, 2019: Generative Modeling by Estimating Gradients of the Data
- 638 Distribution. arXiv preprint arXiv:1907.05600.
- Vannitsem, S., J. B. Bremnes, J. Demaeyer, G. Evans, J. Flowerdew, S. Hemri, S. Lerch,
- N. Roberts, S. Theis, A. Atencia, Z. B. Bouallègue J. Bhend, M. Dabernig, L. D. Cruz, L.
- Hieta, O. Mestre, L. Moret, I. O. Plenkovic, M. Schmeits, J. Ylhäisi, 2021: Statistical
- 642 Postprocessing for Weather Forecasts: Review, Challenges, and Avenues in a Big Data
- 643 World. Bull Amer. Meteor. Soc., **102**, E681–E699.
- Vaswani, A., N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I.
- 645 Polosukhin 2017: Attention Is All You Need, arXiv preprint arXiv:1706.03762.
- Vislocky, R. L., and J. M. Fritsch, 1997: Performance of an advanced MOS system in the
- 1996–97 National Collegiate Weather Forecasting Contest. Bull. Amer. Meteor. Soc.,
- 648 **78**, 2851–2857.
- 649 WMO, 2013: Cascading Process to Improve Forecasting and Warning Services. Bulletin
- 650 *n°*, **62**, 11–15. [Available at https://public.wmo.int/en/resources/bulletin/cascading-
- 651 process-improve-forecasting-and-warning-services.]
- 652

| 654 | List of Figures | | | | |
|-----|--|--|--|--|--|
| 655 | Fig. 1 DeepMedcast framework for training and inference. (a) During the training phase, | | | | |
| 656 | two forecast lead times from the same NWP model (NWP1 at FT = t - Δ t and FT = t + Δ | | | | |
| 657 | are used as input, and the output from the DNN is compared with the same NWP | | | | |
| 658 | model's forecast at FT = t as the ground truth to train the network. (b) During the | | | | |
| 659 | inference phase, predictions from two different NWP models (NWP1 and NWP2) at the | | | | |
| 660 | same lead time (FT = t) are used as input to generate an intermediate forecast between | | | | |
| 661 | the two models at FT = t. | | | | |
| 662 | | | | | |
| 663 | Fig. 2 The recursive application of DeepMedcast, where intermediate forecasts are first | | | | |
| 664 | generated between two NWP models (NWP1 and NWP2, NWP3 and NWP4), followed | | | | |
| 665 | by the creation of an additional intermediate forecast between the outputs of the first two | | | | |
| 666 | pairs. | | | | |
| 667 | | | | | |
| 668 | Fig. 3 The target grid domain for this study. 121 × 151 grids with 0.25-degree × 0.2- | | | | |
| 669 | degree resolution around Japan. The dots on the map represent these grid points. The | | | | |
| 670 | numbers and region names indicated in the figure are used in the case studies in | | | | |
| 671 | Section 3. | | | | |
| 672 | | | | | |
| 673 | Fig. 4 The DNN architecture used in DeepMedcast. The model takes two input channels | | | | |

| 674 | and outputs a single channel. Input data is normalized using the maximum and minimun | | | | |
|-----|---|--|--|--|--|
| 675 | values, and during inference, the same values are applied for the denormalization | | | | |
| 676 | process. | | | | |
| 677 | | | | | |
| 678 | Fig. 5 Comparison of Typhoon LAN predictions by (a) GSM, (b) MSM, (c) the arithmetic | | | | |
| 679 | mean, and (d) DeepMedcast. The forecasts are based on the initial time of 12 UTC on | | | | |
| 680 | 12 August 2023 with a forecast lead time of 51 hours. The black contours indicate mea | | | | |
| 681 | sea-level pressure and wind barbs (units in kt) show surface winds. | | | | |
| 682 | | | | | |
| 683 | Fig. 6 Comparison of DeepMedcast outputs with different input orders for the case in Fig. | | | | |
| 684 | 5. (a) Result when GSM and MSM are provided in that order (same as Fig. 5d). (b) | | | | |
| 685 | Result when the input order is reversed (MSM-GSM). While slight differences are | | | | |
| 686 | present due to network asymmetry, the outputs remain qualitatively identical. | | | | |
| 687 | | | | | |
| 688 | Fig. 7 Comparison of predicted front positions by (a) GSM, (b) MSM, (c) the arithmetic | | | | |
| 689 | mean, and (d) DeepMedcast. The forecasts are based on the initial time of 00 UTC on | | | | |
| 690 | 17 June 2024 with a forecast lead time of 30 hours. The black contours indicate mean | | | | |
| 691 | sea-level pressure, the red contours represent surface temperature, and wind barbs | | | | |
| 692 | (units in kt) show surface winds. The blue, green, purple, and brown dashed lines | | | | |
| 693 | represent the predicted front positions by GSM, MSM, the arithmetic mean, and | | | | |

⁶⁹⁴ DeepMedcast, respectively. Blue circles indicate the location of Tokyo.

| 696 | Fig. 8 Comparison of predicted low-pressure systems by (a) GSM, (b) MSM, (c) the | | | | |
|-----|---|--|--|--|--|
| 697 | arithmetic mean, and (d) DeepMedcast. The forecasts are based on the initial time of 12 | | | | |
| 698 | UTC on 28 August 2024 with a forecast lead time of 75 hours. The black contours | | | | |
| 699 | indicate mean sea-level pressure, the wind barbs (units in kt) represent surface winds, | | | | |
| 700 | the red contours show 850 hPa temperature, and the shaded regions in green and | | | | |
| 701 | yellow highlight areas where the dew-point depression at 850 hPa is below 3° C and | | | | |
| 702 | above 15°C, respectively. | | | | |
| 703 | | | | | |
| 704 | Fig. 9 Comparison of Typhoon KHANUM predictions by (a) GSM, (b) IFS, (c) GraphCast, | | | | |
| 705 | (d) Pangu-Weather, (e) the arithmetic mean, and (f) DeepMedcast. The predictions are | | | | |
| 706 | based on the initial time of 12 UTC on 2 August 2023 with a forecast lead time of 108 | | | | |
| 707 | hours. The black contours indicate mean sea-level pressure and wind barbs (units in kt) | | | | |
| 708 | show surface winds. | | | | |
| 709 | | | | | |
| 710 | Fig. 10 DeepMedcast forecasts based on the initial time of 12 UTC on 2 August 2023 | | | | |
| 711 | with a forecast lead time of 108 hours for Typhoon KHANUN using different orders of | | | | |
| 712 | intermediate forecast generation from the four NWP models: GSM, IFS, GraphCast, and | | | | |
| 713 | Pangu-Weather. (a) Intermediate forecasts are first taken between GSM and IFS, and | | | | |
| | 35 | | | | |

| 714 | between GraphCast and Pangu-Weather, then combined. (b) GSM–GraphCast and | | | | |
|-----|--|--|--|--|--|
| 715 | IFS–Pangu-Weather. (c) GSM–Pangu-Weather and IFS–GraphCast. | | | | |
| 716 | | | | | |
| 717 | Fig. 11 Root mean square error (RMSE) of (a) surface wind speed and (b) surface wind | | | | |
| 718 | direction for DeepMedcast (red), input GSM (blue), and input MSM (green) forecasts, | | | | |
| 719 | verified against AMeDAS observations. The verification period spans one year, from | | | | |
| 720 | January to December 2023, with forecasts initialized at 00, 06, 12, and 18 UTC. | | | | |
| 721 | | | | | |



724

Fig. 1 DeepMedcast framework for training and inference. (a) During the training phase, 725

726

are used as input, and the output from the DNN is compared with the same NWP 727

728 model's forecast at FT = t as the ground truth to train the network. (b) During the

inference phase, predictions from two different NWP models (NWP1 and NWP2) at the 729

same lead time (FT = t) are used as input to generate an intermediate forecast between 730

the two models at FT = t. 731

732



733

The recursive application of DeepMedcast, where intermediate forecasts are first 734Fig. 2 generated between two NWP models (NWP1 and NWP2, NWP3 and NWP4), followed 735by the creation of an additional intermediate forecast between the outputs of the first two 736

pairs.



738



degree resolution around Japan. The dots on the map represent these grid points. The

numbers and region names indicated in the figure are used in the case studies in

742 Section 3.



Fig. 4 The DNN architecture used in DeepMedcast. The model takes two input channels

and outputs a single channel. Input data is normalized using the maximum and minimum

values, and during inference, the same values are applied for the denormalization

process.





Fig. 5 Comparison of Typhoon LAN predictions by (a) GSM, (b) MSM, (c) the arithmetic mean, and (d) DeepMedcast. The forecasts are based on the initial time of 12 UTC on 12 August 2023 with a forecast lead time of 51 hours. The black contours indicate mean sea-level pressure and wind barbs (units in kt) show surface winds.



- Fig. 6 Comparison of DeepMedcast outputs with different input orders for the case in Fig.
- 5. (a) Result when GSM and MSM are provided in that order (same as Fig. 5d). (b)
- 759 Result when the input order is reversed (MSM-GSM). While slight differences are
- present due to network asymmetry, the outputs remain qualitatively identical.



Fig. 7 Comparison of predicted front positions by (a) GSM, (b) MSM, (c) the arithmetic mean, and (d) DeepMedcast. The forecasts are based on the initial time of 00 UTC on 17 June 2024 with a forecast lead time of 30 hours. The black contours indicate mean sea-level pressure, the red contours represent surface temperature, and wind barbs (units in kt) show surface winds. The blue, green, purple, and brown dashed lines

- represent the predicted front positions by GSM, MSM, the arithmetic mean, and
- 769 DeepMedcast, respectively. Blue circles indicate the location of Tokyo.









Fig. 9 Comparison of Typhoon KHANUM predictions by (a) GSM, (b) IFS, (c) GraphCast,
(d) Pangu-Weather, (e) the arithmetic mean, and (f) DeepMedcast. The predictions are
based on the initial time of 12 UTC on 2 August 2023 with a forecast lead time of 108
hours. The black contours indicate mean sea-level pressure and wind barbs (units in kt)
show surface winds.



| 789 | Fig. 10 DeepMedcast forecasts based on the initial time of 12 UTC on 2 August 2023 |
|-----|---|
| 790 | with a forecast lead time of 108 hours for Typhoon KHANUN using different orders of |
| 791 | intermediate forecast generation from the four NWP models: GSM, IFS, GraphCast, and |
| 792 | Pangu-Weather. (a) Intermediate forecasts are first taken between GSM and IFS, and |
| 793 | between GraphCast and Pangu-Weather, then combined. (b) GSM–GraphCast and |
| 794 | IFS–Pangu-Weather. (c) GSM–Pangu-Weather and IFS–GraphCast. |

796



Fig. 11 Root mean square error (RMSE) of (a) surface wind speed and (b) surface wind direction for DeepMedcast (red), input GSM (blue), and input MSM (green) forecasts, verified against AMeDAS observations. The verification period spans one year, from January to December 2023, with forecasts initialized at 00, 06, 12, and 18 UTC. List of Tables

Table 1 Comparison of Typhoon KHANUN predictions from four NWP models, their arithmetic mean, and DeepMedcast. The table presents the predicted central position, central pressure, and maximum wind speed based on the initial time of 12 UTC on 2
 August 2023 with a forecast lead time of 108 hours.

807

Table 1 Comparison of Typhoon KHANUN predictions from four NWP models, their arithmetic mean, and DeepMedcast. The table presents the predicted central position, central pressure, and maximum wind speed based on the initial time of 12 UTC on 2 August 2023 with a forecast lead time of 108 hours.

| | Model | Central position | Central pressure | Maximum wind speed |
|-----|-----------------|------------------|------------------|--------------------|
| (a) | GSM | 31.4°N, 131.1°E | 938 hPa | 68 kt |
| (b) | IFS | 28.9°N, 133.0°E | 966 hPa | 51 kt |
| (C) | GraphCast | 28.0°N, 131.2°E | 977 hPa | 36 kt |
| (d) | Pangu-Weather | 28.8°N, 130.7°E | 975 hPa | 39 kt |
| (e) | arithmetic mean | | 979 hPa | 36 kt |
| (f) | DeepMedcast | 29.3°N, 131.4°E | 964 hPa | 42 kt |

812