

AIFS Single 1.1.0: An update to ECMWF’s machine-learned weather forecast model AIFS

Gabriel Moldovan^{*1}, Ewan Pinnington^{*1}, Ana Prieto Nemesio^{*2}, Simon Lang¹, Zied Ben Bouallègue¹, Jesper Dramsch², Mihai Alexe², Mario Santa Cruz¹, Sara Hahner², Harrison Cook¹, Helen Theissen¹, Mariana Clare², Cathal O’Brien², Jan Polster², Linus Magnusson¹, Gert Mertes¹, Florian Pinault², Baudouin Raoult¹, Patricia de Rosnay¹, Richard Forbes¹, and Matthew Chantry¹

¹European Centre for Medium-Range Weather Forecasts, Shinfield Park, Reading, RG2 9AX, United Kingdom

²European Centre for Medium-Range Weather Forecasts, Robert-Schuman-Platz 3, 53175 Bonn, Germany

Correspondence: Gabriel Moldovan, gabriel.moldovan@ecmwf.int

September 2025

Abstract

We present version 1.1.0 of ECMWF’s Artificial Intelligence Forecasting System (AIFS Single), operational since 25 February 2025. The revised system introduces a bounding-layer framework that enforces physical constraints, such as non-negativity and internal consistency within precipitation and cloud cover variables, alongside expanded training data, revised loss weighting, and an extended set of surface and atmospheric variables. Overall skill improves by 4–6% in the upper air and near-surface variables without degradation of spatial variability. A controlled comparison shows that training data expansion is the dominant source of upper-air skill gains, highlighting the importance of frequent model updates. The bounding framework delivers the largest precipitation improvements, up to 12% and an approximately one-day advantage using a categorical measure of skill. We further show that enforcing precipitation non-negativity resolves a gradient ambiguity at the zero-precipitation boundary under MSE training, explaining the reduction in drizzle bias and the improvements in precipitation.

1 Introduction

Machine-learned weather forecast models have started to rival or outperform physics-based numerical weather prediction (NWP) models in recent years (Pathak et al., 2022; Keisler, 2022; Lam et al., 2023; Chen et al., 2023; Bi et al., 2023; Lang et al., 2024a). For both training

*equal contribution

31 and forecasting, these machine-learned forecast models mostly depend on the Copernicus ERA5
32 reanalysis dataset produced by ECMWF (Hersbach et al., 2020) and operational analysis by
33 ECMWF’s physics-based integrated forecasting system (IFS).

34 ECMWF has developed the artificial intelligence forecasting system (AIFS) (Lang et al.,
35 2024a), its own machine-learned forecast model. After a successful pre-operational test phase
36 running four times daily since October 2023, with forecasts publicly available under ECMWF’s
37 open data policy, AIFS has now transitioned to operational status. The first operational ver-
38 sion, AIFS 1.0.0 replacing AIFS 0.2.1, was implemented on 25 February 2025. The current
39 operational version, AIFS 1.1.0 described here, was released on 27 August 2025 to correct a
40 precipitation forecast issue in the initial version. The model is trained with a mean-squared
41 error (MSE) loss function and is referred to as AIFS Single, to distinguish it from the proba-
42 bilistically trained version, the AIFS ENS (Lang et al., 2024b).

43 Although such MSE-trained forecast models have been shown to smooth forecast fields at
44 longer lead times to avoid the double-penalty of incorrectly positioned weather phenomena
45 (Lam et al., 2023; Ben Bouallègue et al., 2024; Lang et al., 2024a; Bonavita, 2024; Brenowitz
46 et al., 2025), they still display physically robust characteristics (Hakim and Masanam, 2024) and
47 are able to make useful predictions of extreme events (Ben Bouallègue et al., 2024). The cheaper
48 training costs associated with MSE-trained models (compared to probabilistically trained mod-
49 els) make them attractive for prototyping new features and model components.

50 To date, most machine-learned weather forecast models only include a limited subset of
51 forecast variables available from current NWP systems. Here, we include for the first time in
52 the AIFS soil moisture, soil temperature and runoff together with energy sector variables such
53 as cloud cover, 100 metre winds and solar radiation. The choice of additional variables has been
54 guided by utility to users and with considerations of future applications of the model, alongside
55 pragmatic considerations on data availability and readiness. Surface solar radiation and 100-
56 metre wind speeds have been included, important for renewable energy sectors. We added an
57 initial characterization of the land surface with prognostic soil moisture and soil temperature,
58 important for drought forecasting. We also include snowfall, improving the representation of
59 distinct precipitation types in the model. Finally, we have added run-off as a diagnostic model
60 output, pushing towards a hydrological component for the AIFS.

61 Despite their ability to produce skilful forecasts, machine-learned forecast models are prone
62 to producing outputs that violate known physical relationships and limits (e.g., negative pre-
63 cipitation or mass imbalances). In current applications, including the pre-operational version
64 of AIFS, post-processing of forecasts is commonly applied to remove such physical inconsisten-
65 cies. Instead, we propose an additional final layer of activation functions that bound certain
66 variables within physically meaningful limits and enforce physical constraints between related
67 quantities. This simplifies the learning task by constraining the model output space to physi-
68 cally plausible regimes. This bounding strategy also proves particularly beneficial for variables
69 with non-Gaussian distributions, such as precipitation, where the model must effectively dis-
70 tinguish between rain and no-rain states. Enforcing precipitation non-negativity resolves a
71 gradient ambiguity at the zero-precipitation boundary under MSE training, greatly reducing
72 drizzle bias and improving forecast skill in the light-precipitation regime.

73 In this paper we begin by outlining the training setup of the model and how this differs from
74 the previous AIFS version. Then we motivate and describe the new bounding strategy to make
75 the model forecast more physically consistent. We demonstrate the improved performance of
76 the revised AIFS version via evaluation results and selected case studies. We conclude by
77 summarizing main results and future work in the discussion and conclusions.

2 Training

The architecture of AIFS follows an encoder-processor-decoder design. Here, encoder and decoder are attention-based graph neural networks, and the processor is a transformer with a sliding window attention (see Lang et al. (2024a) for details).

The model operates on a reduced Gaussian grid, (N320, approximately 0.25° resolution). The processor (or hidden) grid is an O96 octahedral reduced Gaussian grid (Wedi (2014)) with 40,320 grid points, approximately 1° resolution, and consists of 16 processor layers.

AIFS is trained to produce 6-hour forecasts t_{+6h} using past and present atmospheric states at t_{-6h} and t_0 (from ERA5 or ECMWF’s operational analyses at initialization, or from the model forecast itself). Longer lead times are produced auto-regressively by feeding the model’s predictions back as inputs, a process commonly referred to as rollout.

2.1 Training Schedule

The training is divided into two phases. The first is a pre-training phase, where the model learns to predict the atmospheric state 6 hours ahead (t_{+6h}) using ERA5 analysis at t_{-6h} and t_0 . The second phase, rollout fine-tuning, continues from the pre-trained weights and trains the model to forecast auto-regressively up to 72 hours. Here, the model learns to forecast from its own predictions. Unlike the previous AIFS version, where rollout fine-tuning was first performed using ERA5 and then followed by final fine-tuning on ECMWF operational analysis, we directly use operational analysis for the entire fine-tuning stage. This simplifies the training pipeline, reduces computational costs and is associated with improved forecast performance.

Pre-training is performed on ERA5 data covering the years 1979–2022 (compared to 1979–2020 in the previous AIFS version), using a cosine learning rate (LR) schedule, a batch size of 16, and a total of 260,000 training steps. The LR is linearly increased from 0 to 5×10^{-4} during the first 1,000 steps, then annealed to a minimum of 3×10^{-7} . This is followed by rollout fine-tuning on ECMWF operational analysis from 2016 to 2022, also using a cosine LR schedule and batch size of 16, for approximately 7,900 steps (equivalent to one epoch per rollout step). The LR started at 1.28×10^{-5} and is annealed to the same minimum value of 3×10^{-7} . The rollout length is initially set to 6 hours (1 step) and progressively increased by one step per epoch up to 72 hours (12 steps), following the approach of Lam et al. (2023) and Lang et al. (2024a). We used the AdamW optimizer (Loshchilov and Hutter, 2019) with β coefficients of 0.9 and 0.95. Here, the rollout dataset is extended to eight years of operational IFS analysis (2016–2022), compared with only two years (2019–2020) in the previous AIFS version.

2.2 Variables used in training

The variables used in the new AIFS version are listed in Table 1. As in AIFS 0.2.1, the upper atmosphere is represented by geopotential, horizontal wind components, specific humidity, and temperature at 13 pressure levels: 50, 100, 150, 200, 250, 300, 400, 500, 600, 700, 850, 925, and 1000 hPa. Newly introduced variables are marked with *. We have increased the characterization of the land surface in the model by including new prognostic variables of soil moisture at levels 1 and 2 (swvl1 and swvl2), and soil temperature at levels 1 and 2 (stl1 and stl2), important for drought monitoring and forecasting. A notion of hydrology has been included with runoff (ro), forecast as a diagnostic variable. A second set of variables, related to energy forecasting and clouds, adds real value to the model’s utility. These are forecast diagnostically and include the 100-metre wind components (100u and 100v), surface solar and thermal radiation (ssrd and strd), and cloud cover at various levels (tcc, hcc, mcc, lcc). Finally, snowfall (sf) has been added to complement the set of total precipitation-related variables. An illustration of a selection of these variables can be seen in the forecast presented in Figure 1,

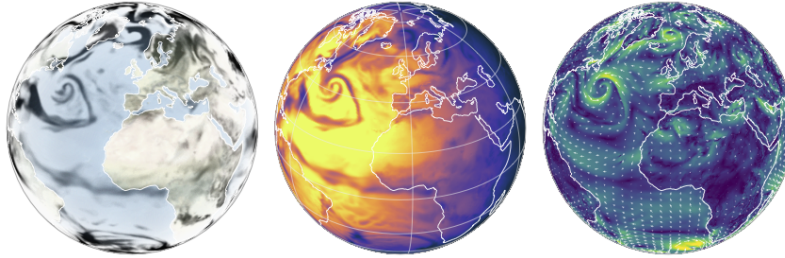


Figure 1: A selection of new variables available from the revised AIFS Single forecasts: cloud cover (left), surface solar radiation (centre), and 100 m wind speed/direction (right). The consistency between these new variables is clear, with areas of higher cloud cover corresponding to lower solar radiation at the surface and consistent weather patterns for 100-metre winds.

124 where the consistency between these new variables is clear, with areas of higher cloud cover
 125 corresponding to lower solar radiation at the surface and consistent weather patterns for 100-
 126 metre winds. These new variables are sourced from the ERA5 reanalysis and IFS operational
 127 data archive, in line with those used in the previous AIFS version (0.2.1).

128 The per variable normalization strategy used in AIFS is summarized in Table 1. Unless
 129 stated otherwise, data is normalized to zero mean and unit variance (z-score normalization). For
 130 some bounded output variables (see Section 3), only standard deviation normalization is applied
 131 to avoid shifting of the absolute zero in the normalized space. The loss function is unchanged
 132 from the previous AIFS version. Table 1 shows the loss scaling factors we use in the revised
 133 AIFS version. Scaling factors were chosen empirically to ensure that all prognostic variables
 134 contribute approximately equally to the loss function, with the exception of vertical velocities
 135 and soil moisture, deliberately down-weighted. Vertical velocity is down-weighted due to known
 136 accuracy limitations in ERA5, particularly in convective regions. Soil moisture receives reduced
 137 weight for similar reasons, and additionally because the transition from ERA5-based pretraining
 138 to operational IFS analysis during fine-tuning introduces distributional inconsistencies; down-
 139 weighting mitigates the influence of this mismatch on training. Furthermore, the loss weights
 140 decrease linearly with height, so that upper atmospheric levels contribute less to the total loss.
 141 The pressure level weights are calculated following $w = \max(\text{pressure level}/1000, 0.2)$, like in
 142 the AIFS-ENS (Lang et al., 2024b). A minimum weight of 0.2 is imposed in the revised version
 143 to avoid assigning excessively low values in the stratosphere.

144 AIFS is trained using data parallelism with a batch size of 16, while each model instance is
 145 distributed across four GPUs within a single node (Lang et al., 2024a). Training was conducted
 146 on the European supercomputer Leonardo (EuroHPC), hosted and managed by Cineca, on
 147 64GB A100 GPUs. Mixed-precision training is used (Micikevicius et al. (2018)), and the full
 148 process takes approximately three days. A 10-day forecast can be produced in about 2 minutes
 149 and 30 seconds on a single A100 40GB GPU, including data input and output.

150 3 Enforcing Model Constraints

151 Machine-learned forecast models for numerical weather prediction show very good forecast
 152 skill, yet they are prone to producing outputs that violate known physical laws or expected
 153 statistical consistency. Unlike traditional numerical models, which are governed by equations
 154 ensuring mass conservation, positivity, or energy bounds, machine-learned forecast models lack

Variable name	Short name	Level type Pressure level (50-1000 hPa) or Surface	Variable type: Prognostic, Diagnostic, Forcing	Normalization	Scaling
Geopotential	z	Pl	P	Z-score	12
Horizontal wind components	u, v	Pl	P	Z-score	0.8, 0.5
Specific humidity	q	Pl	P	Std	0.6
Temperature	t	Pl	P	Z-score	6
Surface pressure	sp	S	P	Z-score	10
Mean sea-level pressure	msl	S	P	Z-score	1
Skin temperature	skt	S	P	Z-score	1
2 m temperature	2t	S	P	Z-score	1
2 m dewpoint temperature	2d	S	P	Z-score	0.5
10 m horizontal wind components	10u, 10v	S	P	Z-score	0.5, 0.5
Total column water	tcw	S	P	Std	1
Volumetric soil water level 1 and 2*	swvl1, swvl2	S	P	None	1, 2
Soil temperature level 1 and 2*	stl1, stl2	S	P	None	1, 10
Total precipitation	tp	S	D	Std	0.025
Convective precipitation	cp	S	D	Std (tp)	0.0025
Snowfall*	sf	S	D	Std (tp)	0.025
Total cloud cover*	tcc	S	D	None	0.1
High cloud cover*	hcc	S	D	None	0.1
Medium cloud cover*	mcc	S	D	None	0.1
Low cloud cover*	lcc	S	D	None	0.1
Runoff*	ro	S	D	Std	0.005
Surface solar radiation downwards*	ssrd	S	D	Std	0.05
Surface thermal radiation downwards*	strd	S	D	Z-score	0.1
100 m horizontal wind components*	100u, 100v	S	D	Z-score	0.1, 0.1
Land-sea mask	lsm	S	F	None	
Orography	z	S	F	Max	
Standard deviation of sub-grid orography	sdor	S	F	Max	
Slope of sub-scale orography	slor	S	F	Max	
Insolation	insolation	S	F	None	
Latitude/longitude (cos/sin)	lat/lon	S	F	None	
Time of day/day of year	local time, julian day	S	F	None	

Table 1: Variables used in the training of AIFS, with their short names, level type, variable type, normalization method, and scaling factors. Variables marked with * were newly introduced compared to AIFS v0.2.1.

155 such guarantees by default. As a result, physically implausible outputs, such as negative
 156 precipitation, can emerge. We show that incorporating constraints into the model design to
 157 enforce physical realism improves forecast skill. In this section, we first identify specific issues
 158 in the output of the previous AIFS version related to total precipitation, and then introduce
 159 a simple yet effective method to bound the model outputs using activation functions. The
 160 proposed method is not restricted to total precipitation but can be equally applied to other
 161 variables.

162 3.1 Lack of Physical Realism in Precipitation Forecasts

163 The previous AIFS version suffers from significant drawbacks in forecasting precipitation. Most
 164 notably, the model’s output is not constrained, leading to a frequent occurrence of negative
 165 values. This is illustrated in Figure 2, which compares the 24-hour accumulated total precip-
 166 itation forecasts from the previous AIFS version, the revised version, GraphCast (Lam et al.,
 167 2023), FuXi (Chen et al., 2023) and an IFS (47r3) 24-hour forecast, for the run initialized
 168 on 01/06/2023 at 00:00 UTC and valid at 02/06/2023 00:00 UTC. The previous AIFS model
 169 and GraphCast show spurious negative precipitation values, which are largely corrected in the
 170 revised AIFS. While negative values can be clipped to zero at inference time (as is done in FuXi
 171 in this figure and thus non visible), their presence highlights a lack of physical consistency in
 172 the model.

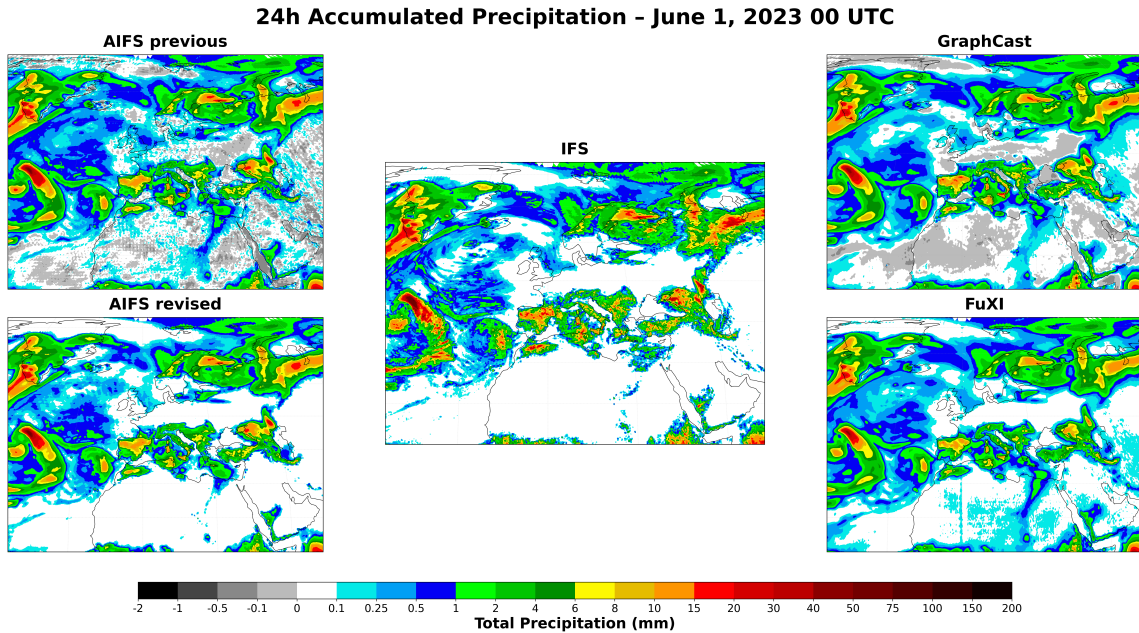


Figure 2: Comparison of 24-hour total precipitation accumulation from five forecasting systems for the forecast issued at 01/06/2023 00:00 UTC and valid at 02/06/2023 00:00 UTC: previous AIFS, revised AIFS, operational IFS, GraphCast, and FuXi. The previous AIFS, GraphCast, and FuXi all exhibit an excess light rainfall, characteristic biases of ML weather models. The revised AIFS, incorporating the bounding layer framework, largely corrects the excess light precipitation issue and provides a precipitation distribution closer to the IFS reference in the light precipitation range.

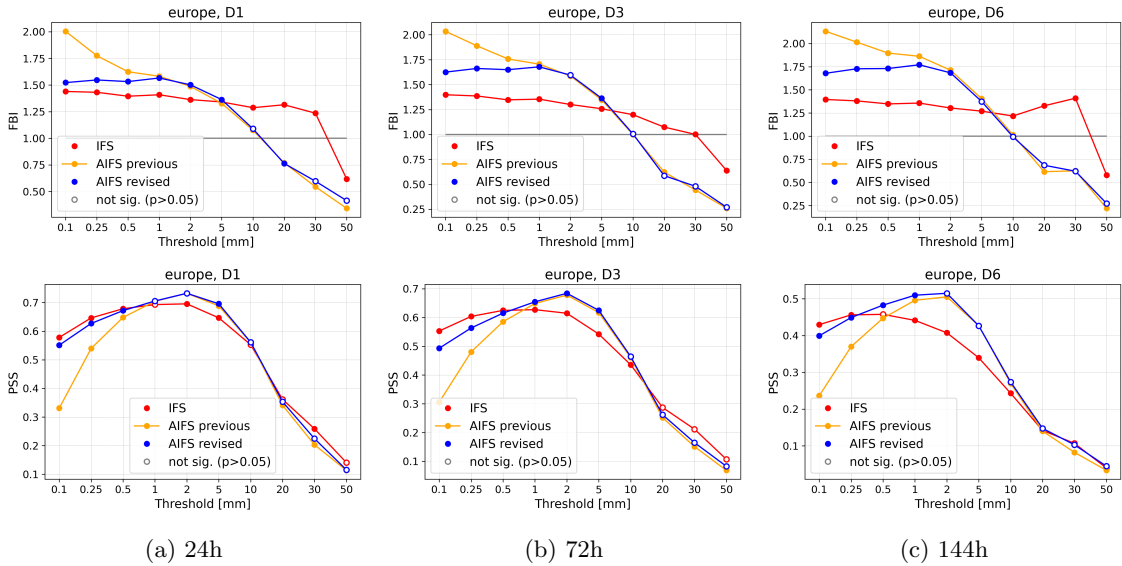


Figure 3: Frequency Bias Index (FBI, top) and Peirce Skill Score (PSS, bottom) for 24-hour accumulated precipitation over Europe as a function of threshold, at forecast day 1, 3, and 6 (left to right). Scores are averaged over all initialisation dates in 2023. Filled markers indicate that the difference relative to the previous AIFS version is statistically significant (paired Wilcoxon signed-rank test, $p < 0.05$); open markers indicate non-significant differences. The previous AIFS version exhibits a pronounced positive frequency bias at low thresholds, consistent with systematic overforecasting of light precipitation.

173 In addition to the negative values, a second noticeable issue, also visible in Figure 2 and
 174 present for all the models but the AIFS revised version, is the excess of light precipitation
 175 in the forecast. The models produce excessive light rain leading to a bias in the forecast.
 176 Similar behaviour has been reported in benchmark studies such as WeatherBench 2 Rasp
 177 et al. (2024), where AI-based systems including GraphCast, Pangu-Weather, and FuXi produce
 178 overly smooth precipitation fields and inflated frequencies of weak events, despite substantial
 179 architectural differences.

180 This is further supported by verification metrics computed against in situ observations
 181 (SYNOP stations). The Frequency Bias Index (FBI) scores for 2023 over Europe (Figure 3)
 182 confirm that the pre-operational AIFS systematically over-forecasts light precipitation events
 183 (< 1 mm). While a similar tendency is present in the IFS, it is considerably more pronounced
 184 in the machine-learned forecast model. At the other end of the distribution, the model tends
 185 to under-forecast more intense precipitation, as indicated by FBI values well below unity for
 186 thresholds exceeding 10 mm. This may be attributed to a well-known characteristic of machine
 187 learning-based forecasts: a tendency to produce overly smooth spatial fields, which can
 188 suppress extremes (Ben Bouallègue et al., 2024; Bonavita, 2024). Additionally, the coarser
 189 native resolution of AIFS (N320 0.25° grid) compared to IFS (0.1° grid) reduces its spatial
 190 representativeness.

191 Convective precipitation forecasts also exhibit similar shortcomings. In addition, there
 192 is a further lack of physical consistency. Convective precipitation represents the part of the
 193 total precipitation that originates from convection, and therefore should always be less than or

194 equal to the total. Figure 4 shows the previous AIFS 24-hour accumulated forecasts of total
 195 and convective precipitation for 02/06/2023. The map displaying the difference between the
 196 two reveals frequent cases in which convective precipitation exceeds total precipitation, which
 197 should not occur.

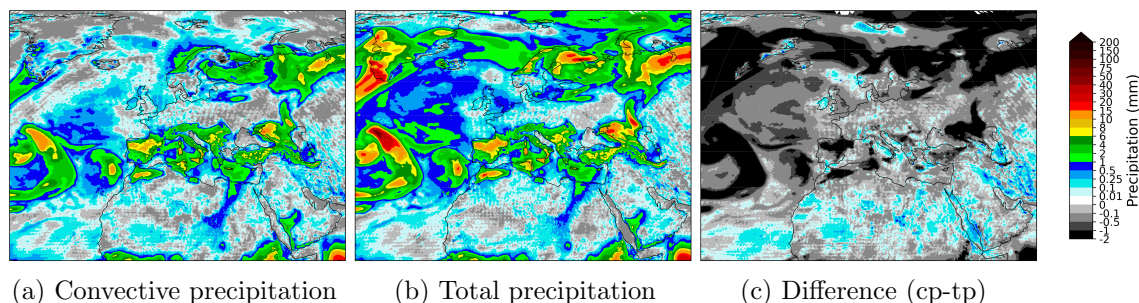


Figure 4: Comparison of 24-hour total and convective precipitation forecast from the previous AIFS version, together with a map showing the difference between the two of them for the forecast issued at 01/06/2023 00:00 UTC and valid at 02/06/2023 00:00 UTC. Positive values (coloured regions) in the difference plot indicate areas where convective precipitation is greater than the total precipitation.

198 The CREDIT platform Schreck et al. (2025) has recently been used to explore physically in-
 199 formed constraints for addressing drizzle bias: Sha et al. (2025b) implemented global mass and
 200 energy conservation schemes as modular constraints within FuXi and demonstrated a direct re-
 201 duction of drizzle bias; a companion study Sha et al. (2025a) further showed that incorporating
 202 terrain-following (hybrid sigma-pressure) can improve extreme precipitation forecasts.

203 Here, we address the drizzle and negative precipitation issue through simplified intervention:
 204 enforcing only the physically admissible output range via a hard-constraint. This approach is
 205 described in Section 3.2. In Section 4.1, we show that this minimal architectural modification
 206 fundamentally reshapes the loss landscape in the vicinity of zero precipitation, eliminating
 207 gradient ambiguity and substantially reducing light-precipitation bias.

208 3.2 Bounding the Outputs with Activation Functions

209 Precipitation has been used as an example to demonstrate the biases present in the forecasts
 210 of some variables. These issues are not only limited to precipitation, but are also observed in
 211 all sparsely distributed variables. This behaviour can be avoided by constraining the output of
 212 the model.

213 There are different strategies one could adopt to enforce physical constraints into the ML
 214 model. More specifically, here we tackled unphysical outputs, and we did not consider other
 215 constraints such as energy or mass conservation. Introducing loss penalties for outputs that
 216 fall outside the known physical bounds can be an effective strategy, and it has the advantage
 217 of not requiring any specific model change. Alternatively, the model could be modified in such
 218 a way as to prevent output from exceeding variable-specific physical bounds. This is usually
 219 referred to as hard-constraining. There are some examples in the literature of hard-constrained
 220 machine-learned models for climate and weather, such as Harder et al. (2024). The authors
 221 apply a softmax function, a generalization of the logistic function, as a hard-constraint for
 222 predicting quantities like atmospheric water content, to enforce the output to be non-negative
 223 in climate downscaling. Other examples can be found in Kent et al. (2025), Bonev et al. (2025)

224 or Subramaniam et al. (2025). Similarly, we argue that hard constraints on the output can be
225 enforced using an activation function.

226 Activation functions can be used in a straightforward way to enforce bounds in the output
227 of machine-learned forecast models. Arguably, the most famous activation function and one
228 we used in this work is the Rectified Linear Unit (ReLU), a nonlinear function defined as:

$$\text{ReLU}(x) = \max(0, x) \tag{1}$$

229 ReLU maps all negative values to zero, effectively enforcing a hard lower bound on the output.
230 For variables requiring both upper and lower bounds, such as concentrations or fractions, the
231 Hard Hyperbolic Tangent (HardTanh) function is an effective choice. It is a piecewise linear
232 approximation of the hyperbolic tangent, defined as:

$$\text{HardTanh}(x) = \begin{cases} 0 & \text{if } x < 0 \\ x & \text{if } 0 \leq x \leq 1 \\ 1 & \text{if } x > 1. \end{cases}$$

233 HardTanh can also be used to enforce consistency between related output variables. For
234 instance, consider the case of convective precipitation (Figure 4), which is predicted indepen-
235 dently of total precipitation in the previous AIFS version. There is a clear relation between
236 the two quantities: convective precipitation is a fraction of total precipitation and should never
237 exceed it. A more physically consistent approach is to map the original convective output
238 to the $[0,1]$ range using a HardTanh layer and to multiply this output by the predicted total
239 precipitation:

$$\text{cp} = \text{HardTanh}(\text{cp}') \times \text{tp}, \tag{2}$$

240 where cp' is the convective precipitation output before the activation layer. This guarantees
241 consistency. This type of constraint, referred to as FractionBounding, is applied to variables
242 related to total precipitation and total cloud cover.

243 Clipping the precipitation output in inference is a possibility and a common practice. This
244 was the case in the pre-operational AIFS model and also reported in other studies, such as
245 Balogh et al. (2024). However, we show that the introduction of bounding in the output during
246 training has benefits beyond simply avoiding slightly negative or unphysical values: it can
247 facilitate the learning of forecasting for sparse and intermittent variables. Bounding effectively
248 decomposes the prediction space into two distinct regions. In the case of total precipitation,
249 the negative space becomes a proxy for forecasting the non-event, while the positive space
250 corresponds to the occurrence of precipitation. This decomposition may, in principle, help the
251 model more easily perform a classification between event and non-event outcomes, a distinction
252 the previous AIFS version struggles with.

253 Table 2 summarises the bounding strategy used in the new version of the AIFS. Since
254 bounding is performed on the normalized space, the choice of the normalization strategy is
255 essential. In particular, variables bounded using a ReLU function were normalized using the
256 standard deviation only, as indicated in Table 1, to avoid offsetting the zero value. Since
257 snowfall and convective precipitation are predicted as fractions of total precipitation, it is
258 necessary to ensure consistent magnitudes in the normalized space. Therefore, cp and sf were
259 scaled using the standard deviation of total precipitation rather than their own. Total cloud
260 cover and soil moisture variables (swvl1 & swvl2) were not normalized, since their range falls
261 within the constraints imposed by the HardTanh bounding ($[0,1]$).

Bounding Type	Range	Variables
ReluBounding	$[0, \infty)$	tp, ro, tcw, ssrd, q(50-1000 hPa)
HardtanhBounding	$[0, 1]$	tcc, swv11, swv12
FractionBounding (w.r.t. tp)	$[0, 1]$	cp, sf
FractionBounding (w.r.t. tcc)	$[0, 1]$	lcc, mcc, hcc

Table 2: Summary of bounding strategies used in the new version of AIFS.

4 Evaluation

Unless otherwise stated, all verification results presented in this section are based on twice-daily forecasts initialised at 00 and 12 UTC for every day of 2023, verified against operational IFS analyses.

The revised AIFS version delivers highly skilled forecasts, as shown by anomaly correlation scores for 2023 in the Northern Hemisphere (Figure 8). In the medium range (3-10 days), AIFS outperforms the IFS by 12 to 24 hours in skill. Forecast skill is also clearly improved compared to the previous AIFS version. This performance gain can be attributed to the combined effect of increased training data, improvements in rollout fine-tuning, the implementation of output bounding, and the inclusion of new prognostic variables. To quantify the specific contribution of expanded training data, we present a controlled comparison in Figure 5. We verify against the operational IFS analysis, which is also used to initialise the forecasts.

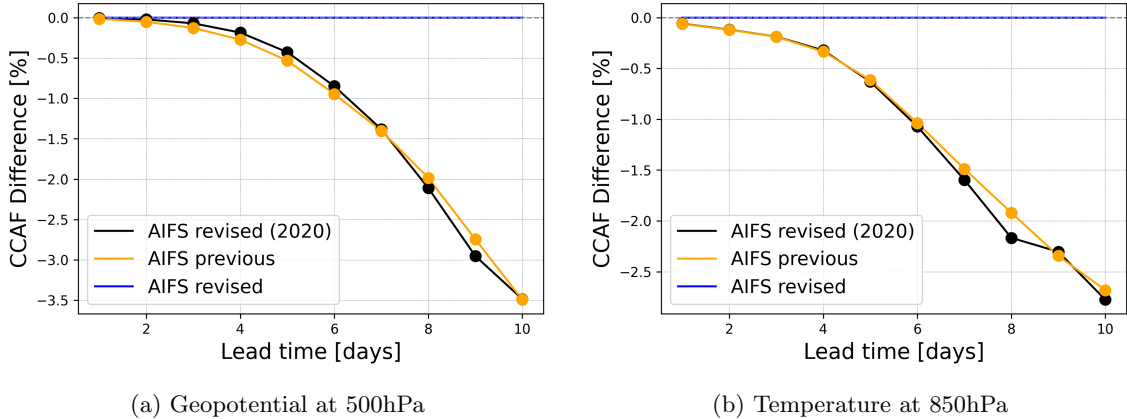


Figure 5: Anomaly correlation skill score difference for Geopotential at 500hPa and Temperature at 850hPa for 2023. This controlled comparison shows: (1) AIFS revised model (full system with all modifications), (2) AIFS revised trained with limited data (ERA5 up to 2020, rollout fine-tuning 2019-2020 only), and (3) AIFS previous version. The close agreement between configurations (2) and (3) demonstrates that the substantial performance gain is primarily attributable to the expanded training dataset (ERA5 1979-2022 and rollout data 2016-2022). Solid points indicate statistically significant differences relative to AIFS revised used as reference (paired Wilcoxon signed-rank test, $p < 0.05$).

As shown in Figure 5, the expanded training dataset contributes to the most important portion of the overall performance gain. This indicates that data availability (ERA5 extended

276 to 2022 and rollout fine-tuning expanded from 2019-2020 to 2016-2022) plays a major role.
277 The remaining improvement stems from other system modifications, including rollout fine-
278 tuning schedules, output bounding layers, and expanded prognostic variables. Due to the
279 high computational cost, a detailed ablation study to isolate the impact of each individual
280 modification beyond data expansion was not performed; thus, the observed improvements
281 represent the cumulative result of these integrated system updates. It should be noted that
282 the close agreement between AIFS revised (2020 data) and AIFS previous in ACC should
283 be interpreted with caution, as these configurations differ in their training protocols: AIFS
284 previous includes a rollout fine-tuning phase on ERA5 which AIFS revised (2020) does not,
285 and uses only 2 years (2019-2020) of operational data for final rollout fine-tuning compared to
286 6 years (2016-2022) in the full revised version. Furthermore, similar ACC scores do not imply
287 equivalent forecast quality. As shown in Figure 6, AIFS revised (2020) exhibits less mesoscale
288 smoothing than AIFS previous despite comparable ACC, indicating that the changes introduced
289 in the revised system do contribute positively to forecast quality in ways not fully captured by
290 ACC alone.

291 Additionally, imposing a minimum on the loss weights in the stratosphere leads to significant
292 improvements in the data-driven forecasts at 100 and 50 hPa (Figure 9). For temperature at
293 100hPa, the new version of the AIFS outperforms the IFS, while for 50hPa wind speed, the gap
294 in skill between the previous version of AIFS and the IFS in the stratosphere is significantly
295 reduced.

296 Forecast skill for key surface variables, such as 2-metre temperature and 10-metre wind
297 speed, verified against SYNOP observations, is similarly improved (Figure 10). Overall, the
298 new AIFS version exhibits improvements of around 4–6 % across all variables, lead times, and
299 pressure levels relative to the previous AIFS version, as shown in the scorecard presented in
300 Figure 7. The performance of the model for tropical cyclone prediction is similar to that of the
301 previous version (see Lang et al. (2024a)), with some small improvements to track position.
302 The training configuration, including a maximum rollout length of 12 (72 hours), was retained
303 from the previous AIFS version, as shown in Section 2.1. This parameter is known to influence
304 spectral characteristics, with longer rollouts leading to enhanced damping. No explicit tuning
305 was performed to target spectral behaviour.

306 The resulting Z500 power spectral density shown in Figure 6 are very similar to those of
307 the previous AIFS across scales, including the 500 km range (zonal wavenumbers 70–90), with
308 slightly improved agreement with the IFS analysis at longer lead times. At the same time, the
309 RMSE-based scorecard (Figure 7) shows overall improvements. Taken together, these results
310 indicate that the skill gains are not achieved at the expense of degraded spatial variability.

311 Figure 11 presents verification metrics for several variables introduced in the new version. In
312 line with those already present in earlier versions, AIFS shows a gain in forecast skill of around
313 one day in the medium range for surface short-wave downwards radiation verified against geo-
314 stationary satellite observation via CMSAF (Pfeifroth et al., 2023) and 100-metre wind speed
315 verified against ECMWF operational analysis, relative to the IFS. The population distribution
316 for total cloud cover verified against SYNOP observations, however, highlights the inherent
317 limitations of MSE-trained AI models. While the observed distribution follows a U-shape,
318 with high frequency at the tails of the distribution (clear skies and overcast conditions), AIFS
319 produces a much flatter distribution, under-predicting these extremes and over-estimating in-
320 termediate values. This behaviour is closely linked to the smoothing effect introduced by the
321 MSE loss function, which tends to penalize large deviations and thereby suppress extremes (see
322 Section 5).

323 The forecasting skill of the model with respect to 24-hour accumulated total precipitation
324 is significantly improved. The new AIFS version is compared against both the previous AIFS
325 version and the operational IFS (cycles 47r3 and 48r1) in Figure 12. The Stable Equitable

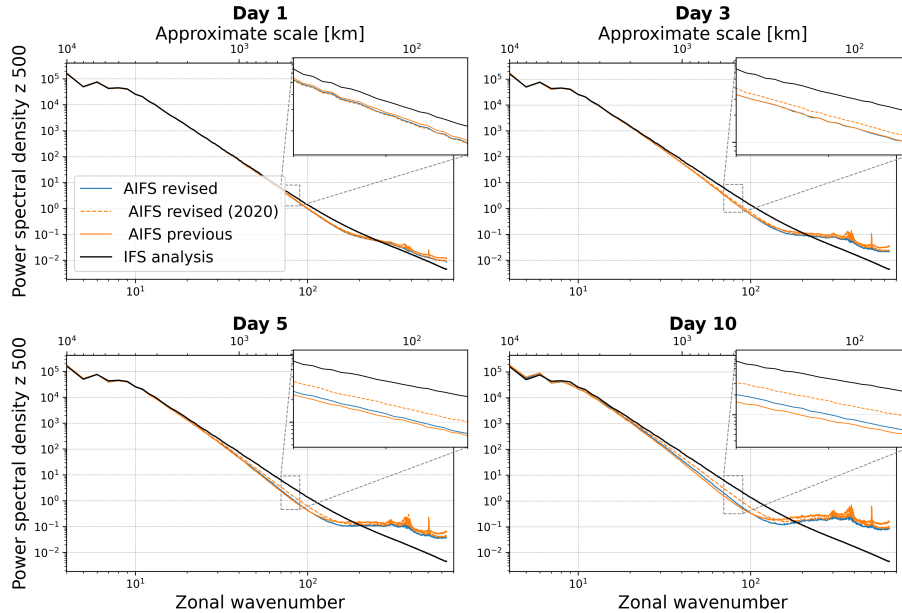


Figure 6: Z500 power spectral density as a function of zonal wavenumber (bottom axis) and approximate horizontal scale in km (top axis) for forecast lead times Day 1, 3, 5, and 10 during JJA 2023. Spectra from the revised AIFS (blue), AIFS revised trained with limited data (ERA5 up to 2020, rollout fine-tuning 2019–2020 only) in dashed orange, and previous AIFS (orange) are compared against the IFS analysis (black). Insets highlight the 450–600 km scale range (zonal wavenumbers 70–90), corresponding to large mesoscale structures. The revised AIFS shows improved agreement with the IFS analysis at large mesoscale structures, particularly at longer lead times, indicating a better representation and retention of mesoscale variance.

326 Error in Probability Space (SEEPS) skill score (Rodwell et al. (2010)) is used as the primary
 327 verification metric, with 24-hour accumulated precipitation SYNOP observations serving as
 328 the reference. Results show a consistent and statistically significant improvement across all
 329 lead times and in the Northern Hemisphere and the Southern Hemisphere. The revised AIFS
 330 demonstrates approximately a one-day gain in forecast skill relative to both IFS and the pre-
 331 vious AIFS version. The forecast fields also exhibit noticeable improvements, as illustrated in
 332 Figure 2. The new version of the AIFS produces no negative values in the output and sub-
 333 stantially reduces light precipitation, aligning more closely with the 24-hour total precipitation
 334 accumulation fields derived from the IFS operational short-range forecasts.

335 Figure 3 reveals where the improvement originates. The Frequency Bias Index (FBI, Wilks
 336 2019), defined as the ratio of predicted to observed event frequency at a given threshold ($FBI = (H+FA)/(H+M)$, where H are hits, FA false alarms, and M misses), and the Peirce Skill Score (PSS, also known as the Hanssen–Kuipers discriminant; Jolliffe and Stephenson 2011), defined as the difference between the probability of detection and the probability of false detection ($PSS = H/(H+M) - FA/(FA+CN)$, where CN are correct negatives), are shown for the Northern Hemisphere for different thresholds. The previous AIFS version exhibits a strong tendency to over-predict light precipitation events (< 1 mm) across all lead times, as shown

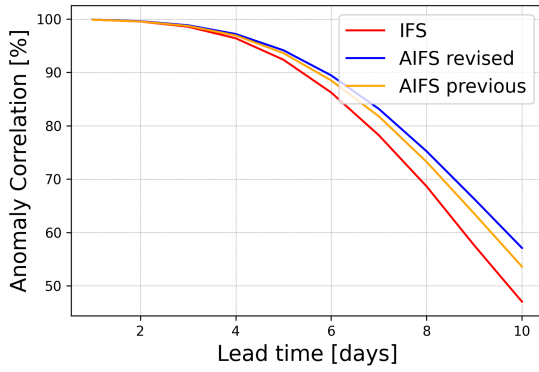


Figure 7: Scorecard comparing forecast scores of AIFS revised versus the previous AIFS version for the whole year of 2023. Forecasts are initialised on 00 and 12 UTC. Relative score changes are shown as function of lead time (day 1 to 10) for northern extra-tropics (n.hem), southern extra-tropics (s.hem) and tropics. Blue colours mark score improvements and red colours score degradations. Purple colours indicate an increased in standard deviation of forecast anomaly, while green colours indicate a reduction. Framed rectangles indicate 95% significance level. Numbers behind variable abbreviations indicate variables on pressure levels (e.g., 500 hPa), and suffix indicates verification against IFS NWP analyses (an) or radiosonde and SYNOP observations (ob). Scores shown are anomaly correlation (ccaf), SEEPS (seeps, for 24h precipitation accumulation), RMSE (rmsef) and standard deviation of forecast anomaly (sdaf).

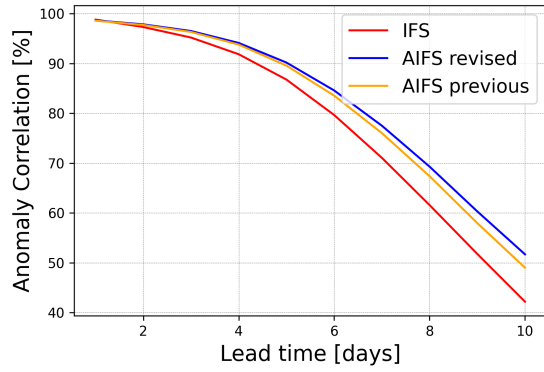
343 by the FBI. This bias is substantially corrected due to the bounding (see Section 4.1) in the
 344 revised AIFS.

345 While the AI model still slightly over-predicts light precipitation compared to the IFS, it
 346 demonstrates competitive skill for light precipitation. The AIFS excels at medium-intensity
 347 events (1–10 mm), with PSS scores significantly higher than those of the IFS. At higher thresh-
 348 olds (> 10mm), corresponding to moderate to heavy precipitation, the AIFS diverges from the
 349 IFS, with a marked under-prediction (FBI < 1). This is likely caused by smoothing introduced
 350 by the loss function, in combination with the model’s coarser spatial resolution.

351 This under-prediction plays an important role in the metrics concerning more extreme
 352 events, since both the previous and the revised AIFS models underperform IFS for thresholds
 353 exceeding 10mm in terms of PSS, but remains competitive. This suggests that although the
 354 AI models predict fewer high-intensity events, their predictions are more accurate when they
 355 do occur. Finally, the revised AIFS shows a marginal improvement in terms of PSS compared

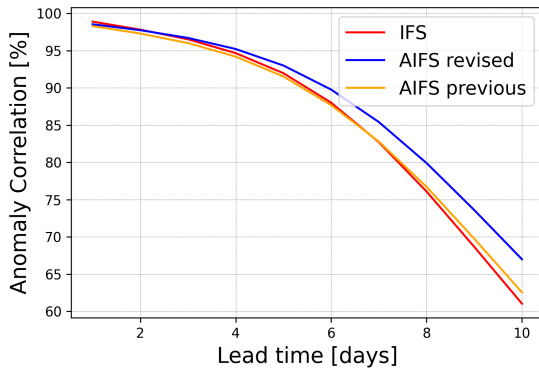


(a) Geopotential at 500hPa

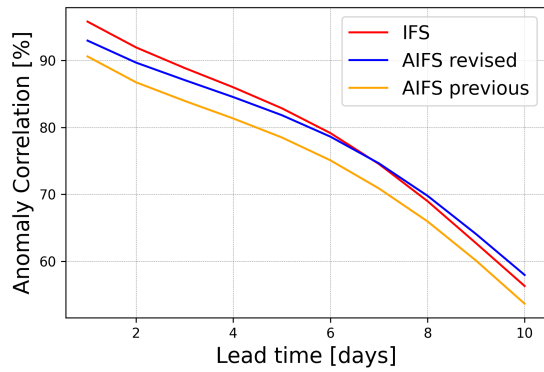


(b) Temperature at 850hPa

Figure 8: Anomaly correlation skill scores for geopotential and temperature at 500hPa and 850hPa, respectively. Skill scores computed for the Northern Hemisphere for the whole of 2023 against IFS analysis. In the medium range, AIFS revised outperforms the IFS by 12 to 24 hours in skill. Forecast skill is also clearly improved compared to the previous AIFS version.



(a) Temperature at 100hPa



(b) Wind Speed at 50hPa

Figure 9: Anomaly correlation skill scores for temperature at 100hPa and wind speed at 50hPa. Skill scores computed for the Northern Hemisphere for the whole of 2023 against IFS analysis. Significant improvements in the revised AIFS forecasts at 100 and 50 hPa when compared against the previous AIFS version.

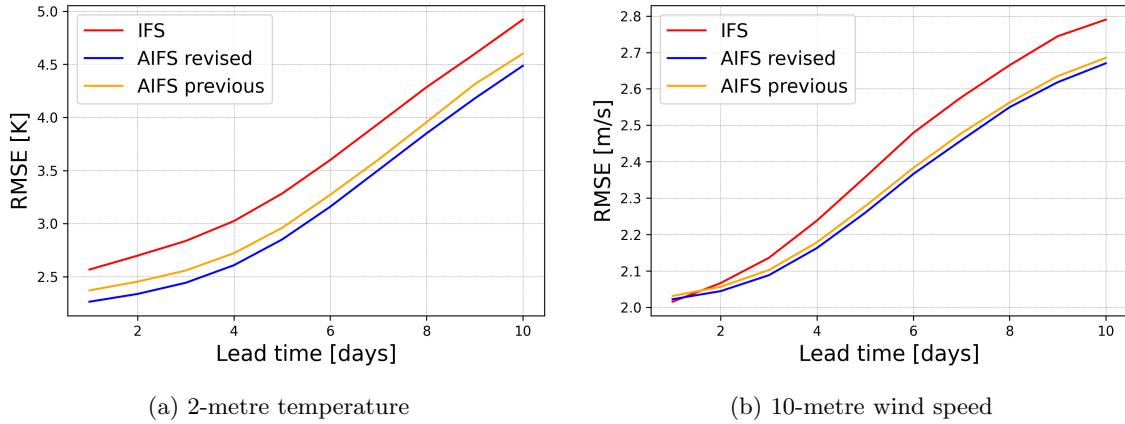


Figure 10: RMSE scores for 2-metre temperature and 10-metre wind speed computed against SYNOP observations over the Northern Hemisphere. The revised AIFS version shows improvement when compared to the previous version of the AIFS.

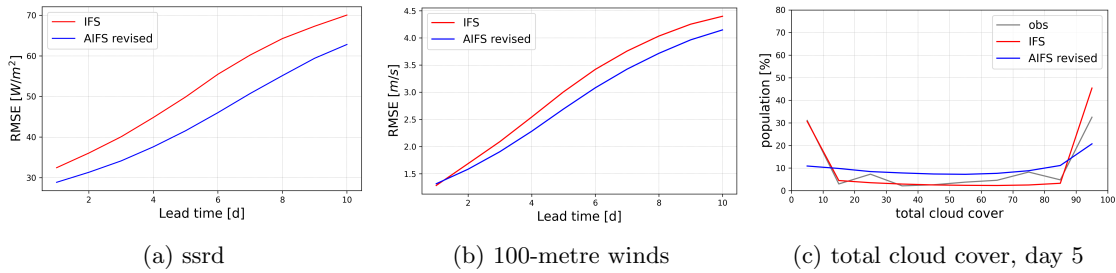


Figure 11: Forecast RMSE computed against operational IFS analysis and distribution comparison for new variables. (a) Surface solar radiation downwards RMSE for March–May (MAM) 2023, (b) 100-metre wind speed RMSE for the full year 2023, (c) Total cloud cover distribution for June–August (JJA) 2023. Blue lines show the AIFS revised and red lines show IFS; observations are shown in grey in panel (c). AIFS shows significant gains in forecast skill in the medium range for surface short-wave downwards radiation and 100-metre winds when compared against the IFS. The mismatch in population distribution for total cloud cover forecast highlights the inherent limitations of MSE-trained AI models.

356
357

against the previous AIFS version, possibly due to improvements in the learning-rate scheduling used for fine-tuning and additional training data.

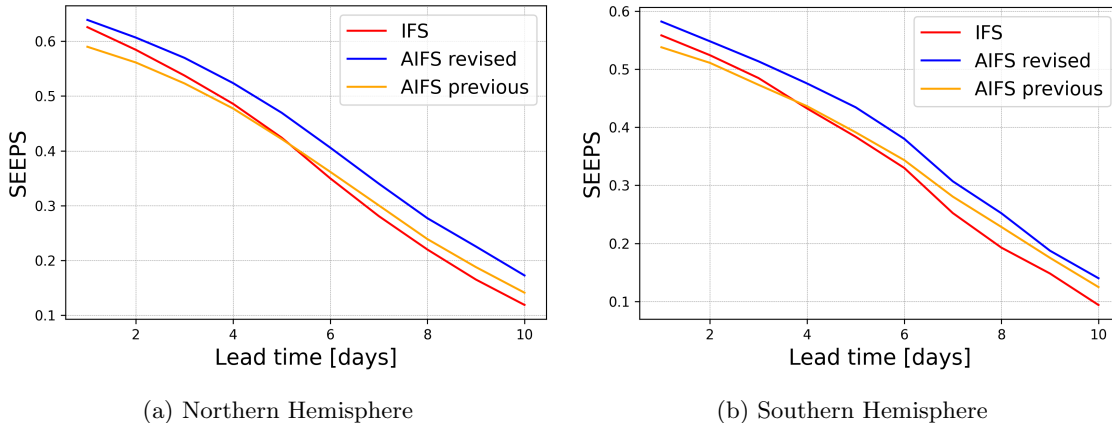


Figure 12: SEEPS skill scores for 2023 based on 24-hour accumulated precipitation from SYNOP observations, comparing the revised AIFS (blue), the previous AIFS version (orange), and the IFS (red) across different regions. Results show a consistent and statistically significant improvement across all lead times and in the Northern Hemisphere and the Southern Hemisphere for the revised AIFS version when compared to the previous AIFS version and the IFS.

358

4.1 Evaluating the effects of bounding on total precipitation

359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380

Overall, the revised AIFS version demonstrates significant improvements in forecasting skill for total precipitation over its predecessor. The bounding of total precipitation transforms the prediction space such that negative values correspond to “no-rain” and positive values to “rain”. This separation enables the model to more effectively distinguish between the two scenarios. It removes the pressure to forecast exactly zero and facilitates the classification task inherent to precipitation forecasting.

Other factors that might improve the precipitation forecast skill in the revised AIFS version are the inclusion of additional variables, the improved learning rate scheduling for rollout fine-tuning and the expansion of the training dataset. To isolate the effect of the bounding mechanism, we retrained the revised AIFS version using the exact same training configuration and data extent, with the sole exception of omitting the bounding layer for total precipitation. This controlled baseline, hereafter referred to as “AIFS revised no-bounding,” allows for a direct comparison between the two models. The SEEPS skill score for the June-July-August 2023 season is shown in Figure 13. The results show that the improvement observed in total precipitation forecast skill in the revised AIFS version can mainly be attributed to constraining the output, since the revised AIFS version without bounding performs similarly to the previous AIFS version.

The physical consistency of convective precipitation forecast in respect to total precipitation can also be evaluated for a given forecast to assess the utility of the FractionBounding strategy used. Figure 14 presents the 24-hour total and convective precipitation accumulation together with a map showing the difference between the two for a forecast issued at 01/06/2023 00:00 UTC and valid at 02/06/2023 00:00 UTC. Unlike the previous AIFS version (Figure 4),

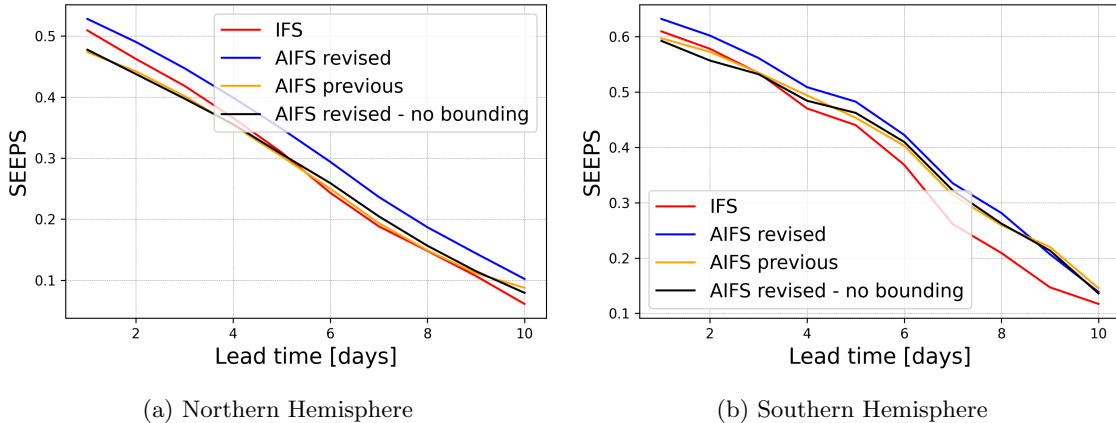


Figure 13: SEEPS skill scores for 2023 JJA comparing revised AIFS (blue), revised AIFS without bounding (black), previous AIFS (orange), and IFS (red) across different regions. The improvement observed in total precipitation forecast skill in the revised AIFS version can mainly be attributed to bounding the output of the model.

381 the convective precipitation forecast is now consistent with the predicted total precipitation
 382 accumulation.

383 To better understand the mechanisms governing total precipitation forecasts in the revised
 384 AIFS configuration, we examine the model’s behaviour in the negative pre-activation space
 385 obtained by removing the final ReLU layer at inference. Figure 15 reveals that this nominally
 386 hidden negative space is neither random nor noisy, but highly structured.

387 At first glance, bounding an output variable with a ReLU activation may appear to intro-
 388 duce a drawback: the negative pre-activation space is not directly penalized, since all negative
 389 values are projected to zero before the loss is evaluated. In principle, changes within this re-
 390 gion do not influence the weight updates. One might therefore expect the negative space to be
 391 uninformative or unstable.

392 Instead, we observe a coherent and physically meaningful organization. Persistently dry
 393 regions, such as the Sahara Desert, exhibit strongly negative pre-activations, while areas
 394 approaching precipitation events transition smoothly toward zero. The model has therefore
 395 learned to encode dryness in the negative space, effectively using it as a latent representation
 396 of the “no-rain” regime.

397 This observation motivates two fundamental questions: (i) why does the negative pre-
 398 activation space contain coherent and physically meaningful structure, and (ii) why does en-
 399 forcing a non-negativity constraint during training improve light-precipitation skill? We argue
 400 that the first arises from the shared latent representation of the atmospheric state learned by
 401 the network, while the second is governed by the symmetry properties of the MSE gradient
 402 near the zero-precipitation boundary.

403 4.1.1 Representation of Dry States in the Negative Space

404 In this study we argue that the structure present in the negative space is an emergent feature
 405 arising from the shared representation of atmospheric states.

406 The model encodes input prognostic (\mathbf{X}_t) and forcing variables (\mathbf{F}_t) into a high-dimensional

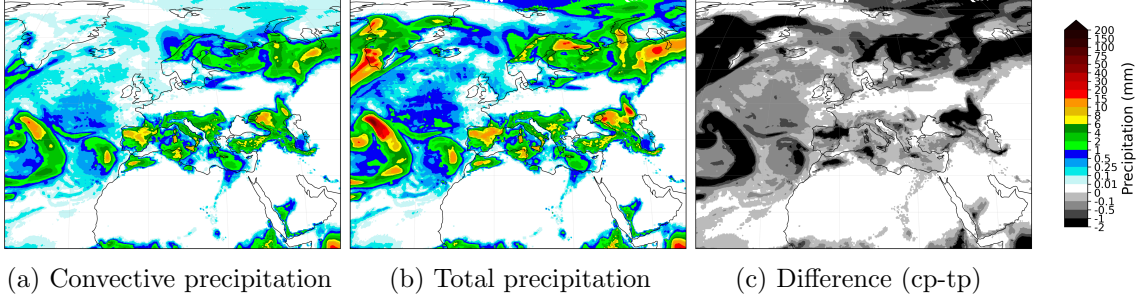


Figure 14: Comparison of 24-hour total and convective precipitation accumulation forecast from the revised AIFS version, together with a map showing the difference between the two of them for the forecast issued at 01/06/2023 00:00 UTC and valid at 02/06/2023 00:00 UTC. Unlike the previous AIFS version (Figure 4), the convective precipitation forecast is now consistent with the predicted total precipitation accumulation and no coloured regions ($cp > tp$) appear in the difference plot.

407 latent space (z_t) via an encoder:

$$z_t = \text{Encoder}(\mathbf{X}_t, \mathbf{F}_t) \quad (3)$$

408 This latent state is evolved to the next time-step through the processor (e.g., via attention-
409 based computations):

$$z_{t+6} = \mathcal{F}(z_t) \quad (4)$$

410 and then decodes back into the physical space to obtain the forecast at t+6 of prognostic
411 (\mathbf{X}_{t+6}) and diagnostic (\mathbf{D}_{t+6}) variables. It is worth mentioning here that z_{t+6} encodes the
412 physical state of all the prognostic variables in a shared representation space and the diagnostic
413 variables are decoded from it. The diagnostic precipitation output is thus produced by a specific
414 decoder head:

$$\eta_{t+6} = \text{Decoder}_{tp}(z_{t+6}) \quad (5)$$

415 where η represents the pre-activation total precipitation. The final physical output is obtained
416 via the bounding layer:

$$tp_{t+6} = \text{ReLU}(\eta_{t+6}) = \max(0, \eta_{t+6}) \quad (6)$$

417 Because Decoder_{tp} maps from a latent space optimized for smooth gradients (z_{t+6}), η in-
418 herits this spatial structure. The precipitation decoder head learns a smooth mapping from
419 the latent space encoding the moisture state of the system to physical precipitation in the
420 positive regime ($\eta > 0$), where gradients are active. Because neural networks are continuous
421 functions biased toward smoothness, this "moisture-to-precipitation" logic naturally extrapo-
422 lates into the negative regime. As moisture variables decrease, the decoder continues to output
423 decreasing values, pushing η into the negative space.

424 While the precipitation head receives no direct gradients when $\eta < 0$, the latent variables
425 that serve as its input are not static. These latent features are shared with prognostic variables
426 (e.g., specific humidity q , total water content tcw , etc) and receive continuous gradient infor-
427 mation from their respective loss functions. Consequently, the negative space of the tp field is
428 "indirectly learned"; it is a projection of a latent space that is being rigorously optimized.

429 Ultimately, this reveals that the optimization of the shared latent space is driven by the col-
430 lective constraints of all output variables. In this framework, the negative pre-activation space

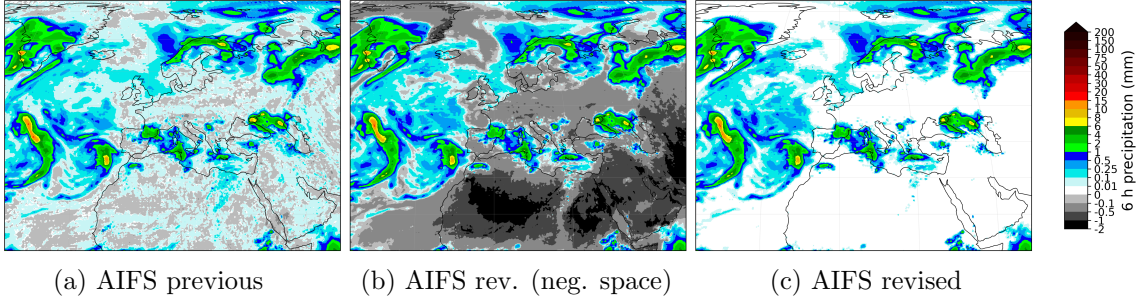


Figure 15: Comparison of 6-hour total precipitation from previous AIFS, revised AIFS without the final ReLU layer to show the negative space, and the standard revised AIFS with the final ReLU layer. Forecasts are initialised at 01/06/2023 00:00 UTC and valid at 01/06/2023 06:00 UTC. Removing the final bounding layer from the AIFS revised model reveals the behaviour of the negative space for the total precipitation variable. The model has implicitly learned to use the negative space as a proxy for “no-rain” classification.

431 for precipitation serves as a “saturation deficit” proxy that is kept physically consistent by the
 432 gradients flowing from prognostic moisture fields. The shared representation of the atmosphere
 433 in the latent space allows the model to maintain a sophisticated, structured representation of
 434 dryness even in the absence of direct precipitation gradients.

435 To provide empirical weight to this mechanistic theory, we investigate the information
 436 content within the pre-activation space η by partitioning the model output into three distinct
 437 physical regimes: the negative (non-precipitating) space, the light precipitation regime (0–0.5
 438 mm/6h), and the moderate precipitation regime (0.5–10 mm/6h).

439 We hypothesize that the pre-activation space η undergoes a fundamental physical decoupling
 440 as it transitions from dry to wet conditions. In the negative (non-precipitating) regime, the
 441 absence of precipitation is a deterministic function of low humidity; thus, the decoder should
 442 preserve a strong linear mapping from the prognostic moisture fields.

443 Conversely, we expect this linear correlation to weaken in the light precipitation regime ($0 <$
 444 $\eta \leq 0.5$ mm). While moisture remains a necessary condition for rain, the exact accumulation
 445 at these low intensities becomes increasingly stochastic, influenced by non-linear factors such
 446 as sub-grid scale turbulence, cloud-base evaporation, and microphysical uncertainties. These
 447 processes act as “interference,” decoupling the surface precipitation from the column moisture
 448 signal.

449 We performed a global correlation analysis on a single forecast issued at 01/06/2023 00:00
 450 UTC. For this experiment, we utilize the AIFS revised model without the final bounding layer
 451 on tp during inference, but activated during training. We focus our analysis on the first 120
 452 hours (5 days) of the forecast.

453 We computed the Pearson correlation coefficient (r) between the pre-activation field η and
 454 five key physical drivers: Total Column Water (TCW), Specific Humidity (q_{1000}), 2m Dewpoint
 455 ($2d$), Mean Sea Level Pressure (MSLP), and mid-tropospheric Vertical Velocity (w_{500}). As
 456 shown in Figure 16, the results reveal a clear regime-dependent physical logic:

- 457 • **Negative Regime ($\eta < 0$):** We observe stable correlations ($r \approx 0.3$) with moisture vari-
 458 ables (q_{1000} , TCW , and $2d$). This confirms that the negative space encodes a structured
 459 representation of the *saturation deficit*, kept physically consistent by gradients flowing
 460 from the prognostic moisture fields.

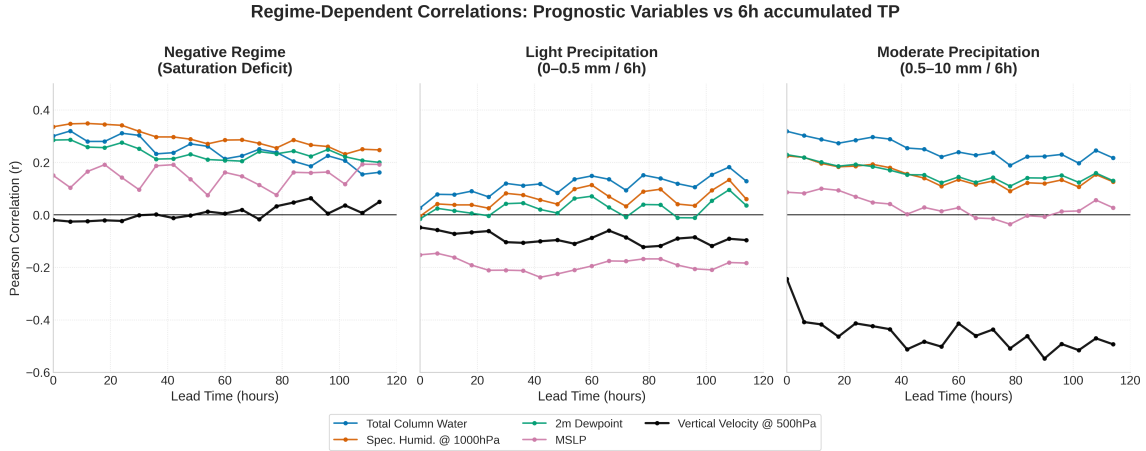


Figure 16: Regime-dependent correlations of pre-activation η (AIFS Revised), for a forecast issued the June 1, 2023 at 00 UTC. Pearson r between η and physical drivers across three regimes: (Left) Negative space ($\eta < 0$): high correlation with moisture variables (q_{1000} , TCW) identifies η as a structured saturation deficit proxy. (Center) Light rain ($0 < \eta \leq 0.5$ mm/6h): systematic weakening of correlation, likely associated with enhanced stochasticity in this regime. (Right) Moderate rain ($1 < \eta \leq 10$ mm/6h): transition to dynamic control, with vertical velocity (w_{500}) as the dominant predictor ($r \approx -0.5$). Analysis covers a 120-hour global forecast.

- 461 • **Light Precipitation** ($0 < \eta \leq 0.5$ mm): Correlation with specific humidity, 2m dewpoint
462 and total column water is substantially reduced in this regime. The weaker relationships
463 are consistent with a lower signal-to-noise ratio and increased sensitivity to
464 small-scale or non-linear processes.
- 465 • **Moderate Precipitation** ($1 < \eta \leq 10$ mm): The model transitions to dynamic control.
466 While moisture correlations remain moderate, Vertical Velocity (w_{500}) emerges as the
467 primary physical driver ($r \approx -0.5$), illustrating the model’s reliance on large-scale ascent
468 to produce deterministic rainfall.

469 While presented as a targeted demonstration of internal model behaviour, the consistency
470 of these signals across lead times suggests that this regime-specific transition is a fundamental
471 structural property of the AIFS architecture. These results demonstrate that the negative
472 pre-activation field encodes valuable information regarding a proxy for saturation deficit. We
473 acknowledge that these correlations are computed from a single 5-day forecast, which limits
474 the temporal sampling. However, the analysis is performed on a Gaussian reduced N320 grid,
475 such that each 6-hourly forecast field contains more than 500,000 spatial evaluation points.
476 Although based on one forecast initialization, the large number of grid-point samples per lead
477 time provides a substantial statistical basis for examining the internal behaviour of the model.

478 4.1.2 Optimization Geometry at the Zero-Precipitation Boundary

479 Having established that the negative pre-activation space encodes physically meaningful information,
480 we now turn to understanding why constraining it during training improves forecast
481 skill for light precipitation. The mechanism can be understood by examining how the Mean

482 Squared Error (MSE) interacts with model outputs in the vicinity of the zero-precipitation
 483 boundary for a non-bounded model:

484 1. **Scenario A (Non-physical negative dry prediction):** The model predicts a non-
 485 physical negative value ($tp = -0.2$ mm) for a dry observation ($tp_{obs} = 0$ mm). The
 486 gradient of the Mean Squared Error (MSE) is:

$$\frac{\partial \mathcal{L}}{\partial tp} = 2(tp - tp_{obs}) = 2(-0.2 - 0) = -0.4 \quad (\text{Push Up}) \quad (7)$$

487 2. **Scenario B (Underprediction):** The truth is light rain ($tp_{obs} = 0.45$ mm), but the
 488 model under-predicts the intensity ($tp = 0.25$ mm). The gradient is:

$$\frac{\partial \mathcal{L}}{\partial tp} = 2(0.25 - 0.45) = -0.4 \quad (\text{Push Up}) \quad (8)$$

489 3. **Scenario C (Overprediction):** The truth is dry or very light rain ($tp_{obs} = 0.05$ mm),
 490 but the model over-predicts the intensity ($tp = 0.25$ mm). The gradient is:

$$\frac{\partial \mathcal{L}}{\partial tp} = 2(0.25 - 0.05) = +0.4 \quad (\text{Push Down}) \quad (9)$$

491 Because non-physical negative dry predictions (Scenario A) and genuine drizzle underpre-
 492 dictions (Scenario B) produce identical upward gradients, the optimizer receives an ambiguous
 493 training signal in the vicinity of zero. The loss provides no information about why the correc-
 494 tion is required — whether it reflects a physical regime transition (dry \rightarrow drizzle) or merely
 495 a violation of the non-negativity constraint. One might expect the model to self-organize by
 496 learning to place dry predictions in a compact negative range — say, around -0.1 mm —
 497 thereby avoiding interference with the light-rain regime. However, this equilibrium is dynami-
 498 cally unstable under MSE. A dry prediction at -0.1 mm receives the same upward gradient as
 499 a genuine drizzle underprediction, so stochastic gradient updates continually push dry samples
 500 toward and across zero. As a result, no stable attractor can form in the negative space.

501 Importantly, the instability is locally asymmetric around $tp = 0$. For small $tp = \epsilon$ with
 502 $|\epsilon| \ll 1$,

$$\frac{\partial \mathcal{L}}{\partial tp} = 2(\epsilon - tp_{obs}).$$

503 In the neighbourhood of zero, the target distribution is one-sided: $tp_{obs} \geq 0$, with strictly
 504 positive drizzle values arbitrarily close to zero but no negative observations. Let

$$\mu = \mathbb{E}[tp_{obs} \mid tp_{obs} \approx 0], \quad \text{with } \mu > 0.$$

505 Then

$$\mathbb{E} \left[\frac{\partial \mathcal{L}}{\partial tp} \right] = 2(\epsilon - \mu).$$

506 Hence the expected gradient is negative for all $\epsilon < \mu$, including the negative space. The only
 507 stationary point of the expected dynamics is $\epsilon = \mu > 0$, which lies strictly on the positive
 508 side. Zero is therefore not a locally stable fixed point under MSE; stochastic gradient updates
 509 induce a systematic drift that transports dry predictions across the boundary into weakly
 510 positive values.

511 As a consequence, dry predictions do not concentrate at a stable negative value but instead
 512 occupy a diffuse region centered on zero, extending into both the negative and weakly positive
 513 ranges. The interval just above zero therefore contains a superposition of displaced dry cases
 514 and genuine drizzle events. This overlap reduces representational separability and compresses
 515 the effective dynamic range available to encode variability within the light-precipitation regime.

516 By enforcing non-negativity through a ReLU constraint during training, negative pre-
 517 activations are projected to zero before loss evaluation. As a result, dry samples no longer
 518 generate corrective gradients within the negative space. Zero becomes a hard boundary rather
 519 than a distributional equilibrium, and the dry regime collapses deterministically onto this
 520 boundary point. The positive axis is therefore freed to encode light-rain variability without
 521 contamination from non-physical corrective gradients.

Bounded vs Non-Bounded AIFS: Output Distribution & Discriminative Capacity in Light Precip

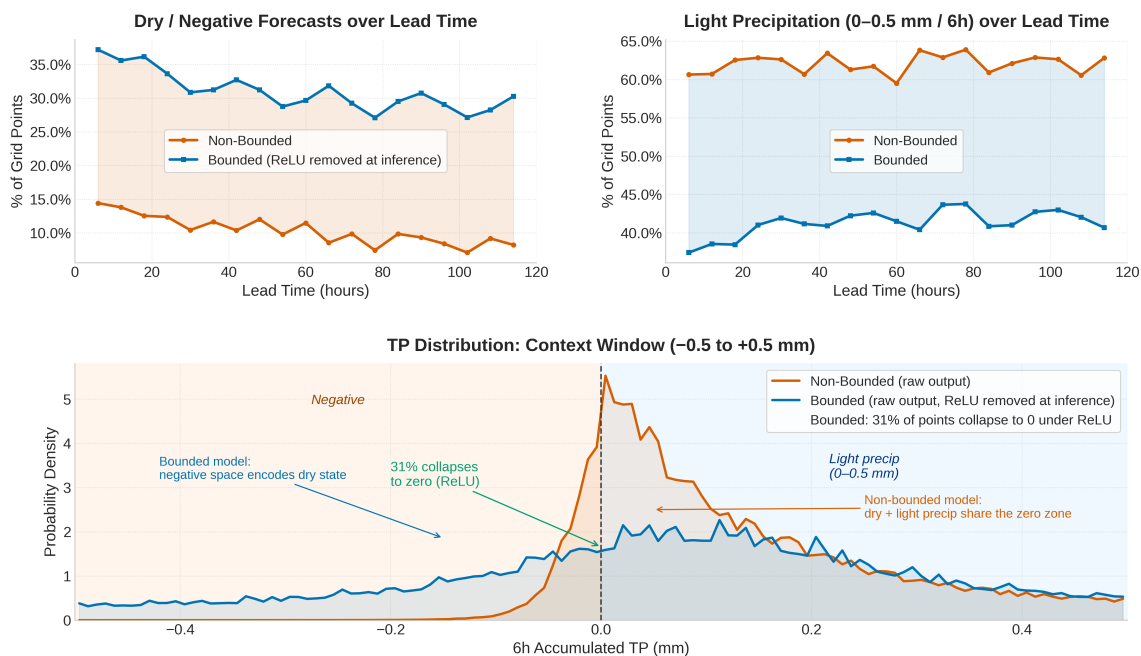


Figure 17: Output distribution and discriminative capacity in the light-precipitation regime for bounded and non-bounded AIFS. The bounded model’s ReLU is removed at inference to expose raw pre-activations. **(Top left)** The non-bounded model produces dry or negative outputs at only $\sim 10\%$ of grid points versus $\sim 30\%$ for the bounded model. A persistent 20-percentage-point gap across all lead times. **(Top right)** The non-bounded model assigns $\sim 60\%$ of grid points to the light-precipitation bin (0-0.5 mm / 6h) versus $\sim 40\%$ for the bounded model, an excess whose magnitude mirrors the dry-detection deficit almost exactly. **(Bottom)** Pre-activation density near zero. The non-bounded model concentrates dry and drizzle cases in an indistinguishable spike around zero; the bounded model distributes dry-state density broadly across the negative space, with 31% of pre-activations collapsing cleanly to zero under ReLU at inference.

522 Figure 17 allows the gradient-ambiguity argument to be verified quantitatively. The three
 523 panels form a closed chain of evidence. The non-bounded model produces dry or negative
 524 outputs at only $\sim 10\%$ of grid points, compared to $\sim 30\%$ for the bounded model. The top-right
 525 panel shows that the non-bounded model’s light-precipitation frequency is inflated by almost
 526 exactly the same ~ 20 percentage points. The non-bounded model is not detecting more drizzle;
 527 it is misclassifying displaced dry events as light rain. The bottom panel reveals the mechanism

528 predicted by the expected-gradient analysis. The non-bounded model produces a narrow spike
529 of density straddling zero, within which the dry and drizzle regimes are superimposed and
530 statistically indistinguishable. The distribution is tightly concentrated near zero but exhibits
531 a slight positive skew, consistent with the theoretical result that the local stationary point of
532 the expected MSE gradient lies at a strictly positive value. In other words, the model attempts
533 to encode dry states in the neighbourhood of zero, yet the systematic upward drift induced by
534 $\mathbb{E}[\partial\mathcal{L}/\partial tp] < 0$ for $tp < \mu$ prevents zero from acting as a stable attractor. The consequence is
535 a persistent displacement of dry samples into weakly positive values, producing the observed
536 excess of light precipitation.

537 Although Figure 17 illustrates a single 5-day forecast, the behavior is systematic rather
538 than case-specific. This interpretation is reinforced by the Frequency Bias Index (FBI) and
539 Peirce Skill Score (PSS) shown in Figure 3 of the main article. The non-bounded configuration
540 exhibits a pronounced positive frequency bias in the light-precipitation category, together with
541 degraded discrimination skill, consistent with systematic misclassification of dry grid points as
542 drizzle.

543 The mechanism described here provides a refined interpretation of recent findings in AI-
544 driven precipitation forecasting. Sha et al. (2025) reported that drizzle bias is substantially
545 reduced when physical constraints are applied, whereas terrain-following coordinates alone do
546 not mitigate drizzle bias but instead improve extreme precipitation forecasts. Notably, their
547 constraint framework combines global conservation principles with an explicit non-negativity
548 correction.

549 The present analysis isolates the role of non-negativity enforcement and demonstrates that it
550 addresses a fundamental gradient asymmetry at the zero-precipitation boundary. This mech-
551 anism operates at the level of local optimization dynamics and provides a distinct, mech-
552 anistically interpretable pathway for drizzle reduction. While Sha et al. (2025a) demonstrate
553 effectiveness of combining non-negativity with global conservation constraints, our analysis
554 suggests that non-negativity merits investigation as an independent design element. The rela-
555 tive contributions of boundary enforcement versus conservation-based regularization, and their
556 potential architecture dependence, remain important questions for future work.

557 4.2 Case Studies

558 Headline verification scores for the revised AIFS show significant improvements over the conven-
559 tional numerical weather prediction model. However, building trust in AI forecasting requires
560 more than strong overall metrics. Forecasters place great importance on the ability of the
561 model to accurately and reliably predict weather phenomena. They also value physically plau-
562 sible outputs and recognizable weather patterns. To support this, we show below selected case
563 studies.

564 4.2.1 Storm Éowyn

565 Storm Éowyn was an unusually strong winter storm and blizzard, initially impacting much of
566 the Gulf Coast of the United States between January 20 and January 22, 2025. This storm broke
567 snowfall records at a number of reporting stations (Thiem and Collins, 2025) and represented
568 an extreme out-of-training-distribution event with no clear analogies in the ERA5 reanalysis
569 or the IFS Operational analysis dataset.

570 Figure 18 shows the AIFS and IFS forecasts at decreasing lead times for the affected area
571 versus the corresponding IFS short-range forecast. The AIFS delivers an accurate forecast of
572 snowfall for this extremely rare event. This showcases the ability of the model to accurately
573 interpret meteorological patterns and forecast physically plausible events, even if they are far

574
575

from the training data. The AIFS predicted the event with a lead time of 10 days, earlier than the IFS.

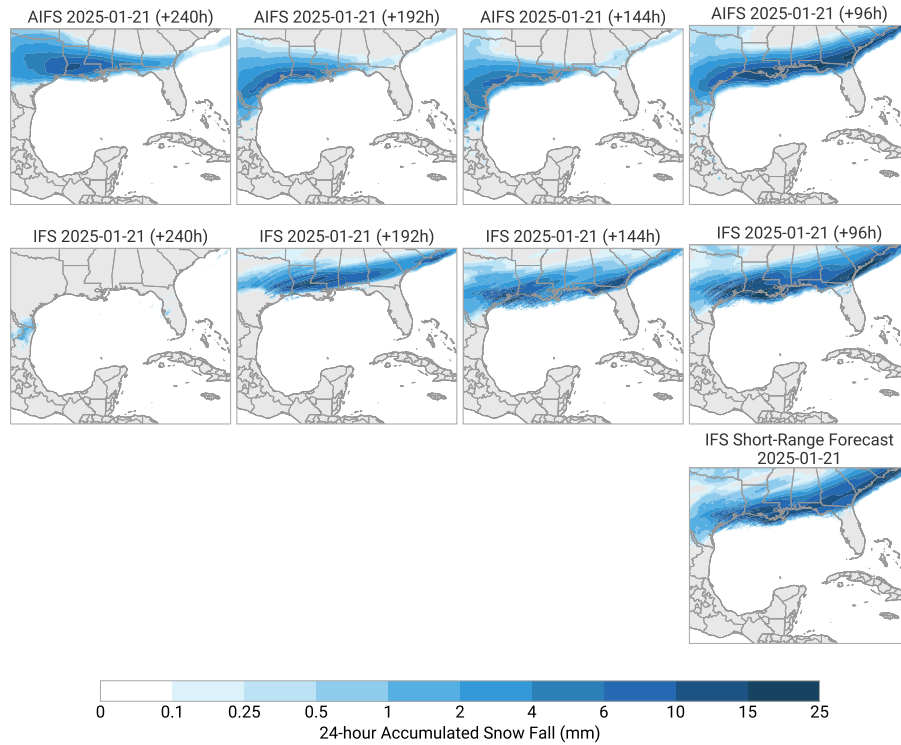


Figure 18: Snowfall forecasts for AIFS (top row) and IFS (middle row) over the Gulf Coast of America at 10, 8, 6 and 4 day lead times from left to right respectively, against IFS short-range forecasts for the snowfall event (bottom row). The figure shows how the snowfall event was forecast accurately four days ahead by both the IFS and AIFS. The AIFS forecasted the event even 10 days ahead.

576
577

4.2.2 Tropical Low and extreme precipitation totals in Queensland Australia

578
579
580
581
582
583
584
585
586
587

Starting in late January 2025, a slow-moving summer storm brought exceptional rainfall along the northeastern coast of Queensland, Australia. Within a week, rainfall accumulation totalled more than 1000 millimetres in some areas, according to the Bureau of Meteorology as reported in NASA Earth Observatory (2025). The city of Townsville saw the equivalent of six months of rain in just three days and the largest weekly rainfall total was measured at a gauge in the Cardwell Range, southwest of Tully, where nearly 1700mm fell (NASA Earth Observatory (2025), Bureau of Meteorology measurements). Figure 19 compares forecasts from AIFS and IFS against the IMERG Huffman et al. (2023) final product for the period 01/02/2025–03/02/2025. Both model forecasts were initialized on 30/01/2025, two days prior to the event. The Cardwell Range is indicated by a black star, and the city of Townsville by a

588 cyan star. Both IFS and AIFS successfully captured the event, with 24-hour rainfall accumu-
 589 lations exceeding 300 mm in some regions. However, the AIFS forecast exhibits a somewhat
 590 persistent signal in the 5-day lead time, predicting very high rainfall totals near the Cardwell
 591 Range. This highlights that, despite AIFS’s tendency toward excessive spatial smoothing, it
 592 remains capable of accurately forecasting extreme events at medium range.

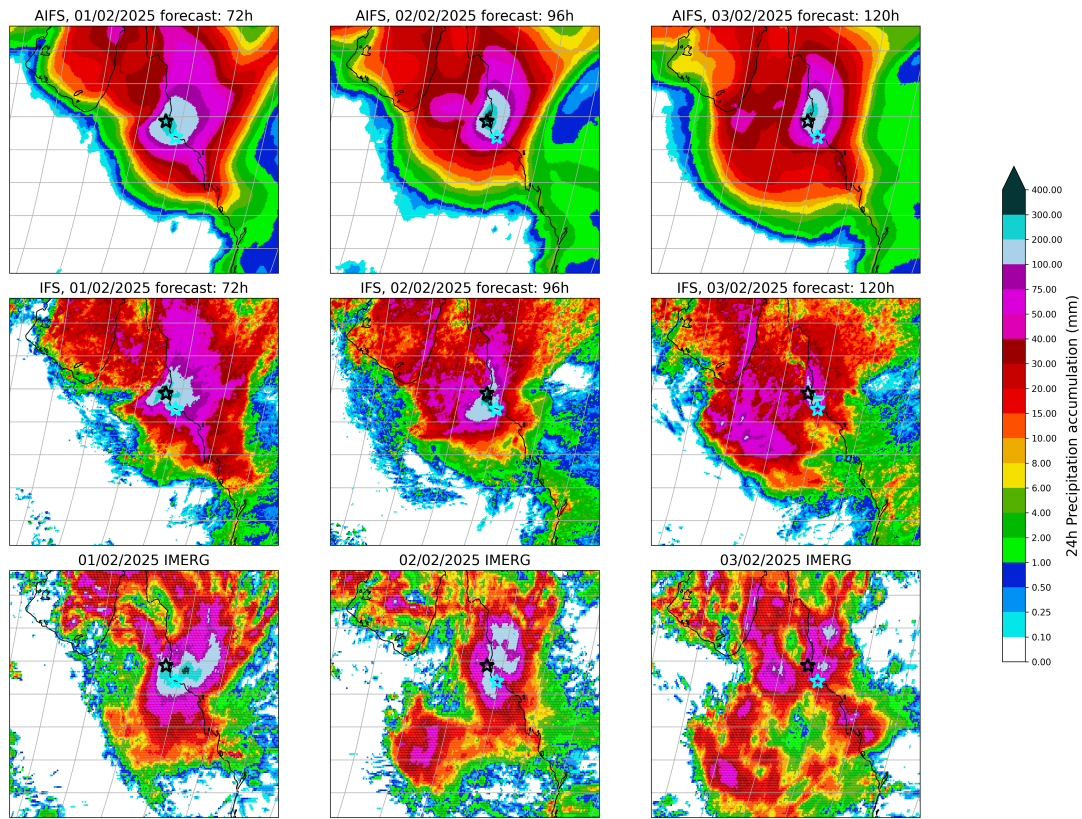


Figure 19: 24-hour accumulated precipitation forecasts from the AIFS (top row) and IFS (middle row) models, compared with IMERG observational data (bottom row) over northeastern Queensland for 01/02/2025 to 03/02/2025. Forecasts are initialised on 30/01/2025. The black star marks the Cardwell Range, where rainfall totals exceeded 1600 mm over the week, and the cyan star marks the city of Townsville. Both models captured the core of the extreme rainfall event, with accumulations exceeding 300 mm in 24 hours in some areas.

593 5 Discussion and conclusion

594 The revised AIFS version (1.1.0) presented here improves upon the pre-operational release
 595 through a revised training regime with more data, new forecast variables, improved strato-
 596 spheric loss weights, and a bounding strategy that enforces physical constraints on the output
 597 variables. Overall, this leads to improvements of around 4–6 % across all variables, lead times,

598 and pressure levels. The largest improvements, up to 12% gains in normalized difference in the
599 short range, are observed in total precipitation forecasting, which benefits from the newly intro-
600 duced bounding. We showed that this has a significant impact on the prediction of no rain and
601 light precipitation. The model displays good forecast performance for out-of-training-sample
602 case studies, accurately capturing extreme precipitation and snowfall events.

603 Data plays a crucial role in the performance of AI models. Most of the improvements
604 non-related to precipitation in the revised version of the AIFS stem from the expansion of
605 the training dataset and the use of more recent operational ECMWF analyses for rollout fine-
606 tuning, as demonstrated by the controlled comparison in Figure 5. Since the AIFS relies on
607 these analyses for real-time forecasting, it is important to fine-tune them regularly using up-
608 to-date data. Regular fine-tuning with recent ECMWF analyses helps the models to adapt to
609 shifts in the data due to new IFS model cycles.

610 Recent global AI forecasting systems, including GraphCast, Pangu-Weather, FuXi, and
611 CREDIT, have reported persistent challenges in representing light precipitation. Positive fre-
612 quency bias in the drizzle regime appears to be a recurring feature across models trained
613 with symmetric regression losses on strictly non-negative, intermittent variables. Although
614 these systems differ substantially in backbone architecture, from graph neural networks to
615 transformer-based designs and modular physically constrained frameworks, the drizzle prob-
616 lem appears largely independent of architecture. Instead, it is closely tied to how precipitation
617 is parameterized and constrained during training. Physical constraint methodologies offer mul-
618 tiple pathways for mitigating precipitation biases. Global conservation schemes may reduce
619 drizzle indirectly by regulating total moisture budgets. The present analysis suggests that non-
620 negativity enforcement addresses a more fundamental issue: the local gradient asymmetry at
621 the zero boundary and the superposition of dry and wet states around zero. By introducing
622 a hard geometric boundary at zero, the optimization landscape is reshaped such that dry and
623 wet regimes become separable. This mechanism operates independently of large-scale conser-
624 vation principles and may therefore represent a structural requirement for stable training of
625 intermittent variables under MSE. Alternative activation functions such as LeakyReLU, which
626 scale negative inputs by a small factor α (typically 0.01), would permit gradient flow in the
627 negative space while heavily attenuating the loss contribution from dry predictions (by a factor
628 of α^2). We expect that similar regime separation would still emerge, since the cost of placing
629 dry states deep in the negative space becomes negligible. The main practical difference is that
630 LeakyReLU produces non-physical slightly negative output values at inference, requiring post-
631 processing clipping. More broadly, alternative formulations that explicitly decouple the dry
632 and wet regimes during training, such as asymmetric loss functions or dedicated classification
633 heads for the no-rain state, represent promising directions for future work.

634 The bounding strategy presented here enforces physical realizability, non-negativity, bound-
635 edness, and inter-variable consistency, but does not impose global conservation of mass or en-
636 ergy. For the medium-range timescales considered in this work (up to 10 days), we expect
637 conservation violations to remain small relative to forecast errors dominated by chaotic er-
638 ror growth, though a systematic quantification of mass and energy drift over extended AIFS
639 integrations remains to be carried out.

640 Rollout fine-tuning emerges as an important factor shaping forecast behaviour, including
641 the degree of spatial smoothing in the outputs. As the model is trained on extended lead times
642 and optimised using a mean squared error objective, some degree of smoothing is expected.
643 Training hyperparameters such as learning rate scheduling, number of optimisation steps, and
644 rollout configuration can influence this behaviour and warrant further systematic investigation.
645 In the present study, the training configuration, including a maximum rollout length of 12, was
646 retained from the previous AIFS version to ensure consistency. The resulting Z500 power spec-
647 tra (Figure 6) are broadly comparable to those of the previous model across scales, including

648 the 500 km range, with slightly improved agreement with the analysis at longer lead times.
649 Importantly, these comparable spectral characteristics are achieved alongside overall improve-
650 ments in RMSE-based skill (Figure 7). This indicates that the skill gains are not obtained at
651 the expense of degraded spatial variability. While more aggressive rollout strategies may fur-
652 ther optimise headline verification scores, understanding their impact on spatial characteristics
653 remains an important area for future work.

654 Alongside making updates to the training schedule, we have also added new variables to
655 the AIFS while achieving improvements in forecast skill for headline atmospheric metrics. In
656 particular, the inclusion of soil moisture and soil temperature as prognostic variables represents
657 an initial step toward a more complete Earth system representation within AIFS. Targeted
658 ablation studies are planned as the land-surface component is extended in future versions.
659 However, it remains to be seen if adding more variables and earth-system components will
660 eventually require an increase to the latent space of the model. The additional earth-system
661 and energy-sector variables in AIFS establish a foundation for future extensions, including
662 ocean and wave components, expanding the number of cryospheric processes with enhanced
663 snow modelling, and increasing the hydrological capabilities of the model. These new variables
664 are currently taken from a consistent data source with the rest of the model variables. In the
665 future, there is the potential to look at datasets tailored to specific earth-system components,
666 such as ERA5-Land (Muñoz Sabater et al., 2021) and the ocean and sea-ice reanalysis system
667 (ORAS6) (Zuo et al., 2024).

668 AIFS currently operates at approximately 0.25° spatial resolution with a 6 hour timestep,
669 and future work will focus on increasing both spatial and temporal resolution.

670 The AIFS development has now transitioned to the new Anemoi framework (Lang et al.,
671 2024a; Nipen et al., 2024; Wijnands et al., 2025). Anemoi provides tools for the whole data-
672 driven modelling workflow, from the generation of training datasets, to scalable probabilistic
673 training (Lang et al., 2024b) and running real-time inference with such models. Anemoi also
674 allows for the cataloguing and archiving of model and data checkpoints to ensure reproducibility
675 and traceability of training and inference runs and ensure that any models developed within
676 this framework have a clear lineage. The Anemoi framework is now being used by an increasing
677 number of Member States of ECMWF and collaborating organisations supported by ECMWF.

678 After a successful experimental phase, AIFS has transitioned to operational status at
679 ECMWF on the 25th of February 2025. It is supported 24/7 alongside ECMWF’s physics-
680 based system, the IFS. The MSE trained model is labeled AIFS Single, and its forecasts are
681 available earlier than the ones from the physics-based model chain, due to the fast runtime
682 of AIFS. Results presented in this paper show that AIFS forecasts are highly skilful and they
683 outperform the IFS forecasts across the vast majority of lead times and variables. They high-
684 light the relevance of AIFS for weather prediction. Future developments will focus on including
685 more surface variables and exploring a wider range of applications such as climate reanalysis.
686 The operational release of the AIFS demonstrates the commitment of ECMWF to pursue the
687 best possible weather forecasts with both physics-based and machine learning methods.

688 6 Code and Model Availability

689 AIFS version 1.1.0 was fully trained using the Anemoi framework [https://github.com/ecmwf/](https://github.com/ecmwf/anemoi)
690 [anemoi](https://github.com/ecmwf/anemoi). The frozen versions of the Anemoi modules used for training, together with the config-
691 uration files and the trained model checkpoint, are available in the permanent archive European
692 Centre for Medium-Range Weather Forecasts (2025) under DOI [https://doi.org/10.5281/](https://doi.org/10.5281/zenodo.17349820)
693 [zenodo.17349820](https://doi.org/10.5281/zenodo.17349820). The model weights for version 1.1.0 are also available on the project page on
694 Hugging Face <https://huggingface.co/ecmwf/aifs-single-1.1> under a Creative Commons

695 Attribution 4.0 International (CC BY 4.0) licence and DOI [https://doi.org/10.57967/hf/](https://doi.org/10.57967/hf/6415)
696 6415 (ECMWF, 2025a).

697 The AIFS Single model operational forecasts are freely available under ECMWF’s Open
698 Data Creative Commons licence (<https://www.ecmwf.int/en/forecasts/datasets/open-data>)
699 and DOI <https://doi.org/10.21957/open-data> (ECMWF, 2025b) and forecast charts can
700 be seen at <https://charts.ecmwf.int/?query=aifs-single>. Further details on the model’s
701 operationalization and data dissemination can be found at [https://confluence.ecmwf.int/](https://confluence.ecmwf.int/display/USS/Implementation+of+AIFS+Single+v1.0)
702 [display/USS/Implementation+of+AIFS+Single+v1.0](https://confluence.ecmwf.int/display/USS/Implementation+of+AIFS+Single+v1.0).

703 Author Contributions

- 704 • **Experiment design and execution:** GMo*, EP*, APN*, SL, MCh
- 705 • **Model evaluation:** GMo*, EP*, APN*, ZBB*, LM, SL, MCh
- 706 • **Framework development (Anemoi):** SL, JD, MCh, MA, APN, MSC, SH, HC, HT,
707 MC, CO, JP, GMe, FP, BR, GMo, EP
- 708 • **Manuscript writing:** GMo, EP, SL, APN with input from all co-authors

709 *Equal Contribution

710 7 Acknowledgements

711 We acknowledge PRACE for awarding us access to Leonardo, CINECA, Italy. We acknowledge
712 the EuroHPC Joint Undertaking for awarding this work access to the EuroHPC supercomputer
713 MN5, hosted by BSC in Barcelona through a EuroHPC JU Special Access call. Ewan Pin-
714 nington’s contribution is funded under the CERISE project (grant agreement No101082139),
715 CERISE is funded by the European Union. Ana Prieto Nemesio’s contribution is partially
716 funded under the RODEO project (grant agreement: No101100651), RODEO is funded by the
717 European Union. Views and opinions expressed are however those of the author(s) only and do
718 not necessarily reflect those of the European Union or the Commission. Neither the European
719 Union nor the granting authority can be held responsible for them.

720 References

- 721 Blanka Balogh, David Saint-Martin, and Olivier Geoffroy. Online test of a neural network deep
722 convection parameterization in arp-gem1, 2024. URL <https://arxiv.org/abs/2410.21920>.
- 723 Zied Ben Bouallègue, Mariana C A Clare, Linus Magnusson, Estibaliz Gascón, Michael Maier-
724 Gerber, Martin Janoušek, Mark Rodwell, Florian Pinault, Jesper S Dramsch, Simon T K
725 Lang, Baudouin Raoult, Florence Rabier, Matthieu Chevallier, Irina Sandu, Peter Dueben,
726 Matthew Chantry, and Florian Pappenberger. The rise of data-driven weather forecasting: A
727 first statistical assessment of machine learning-based weather forecasts in an operational-like
728 context. *Bulletin of the American Meteorological Society*, 2024. ISSN 1520-0477. doi: [doi: doi.](https://doi.org/10.1175/BAMS-D-23-0162.1)
729 [org/10.1175/BAMS-D-23-0162.1](https://doi.org/10.1175/BAMS-D-23-0162.1). URL <http://dx.doi.org/10.1175/BAMS-D-23-0162.1>.
- 730 K. Bi, L. Xie, H. Zhang, et al. Accurate medium-range global weather forecasting with 3D
731 neural networks. *Nature*, 619:533–538, 2023. doi: [10.1038/s41586-023-06185-3](https://doi.org/10.1038/s41586-023-06185-3).

732 Massimo Bonavita. On some limitations of current machine learning weather prediction models.
733 *Geophysical Research Letters*, 51(12):e2023GL107377, 2024. doi: [https://doi.org/10.1029/](https://doi.org/10.1029/2023GL107377)
734 [2023GL107377](https://doi.org/10.1029/2023GL107377).

735 Boris Bonev, Thorsten Kurth, Ankur Mahesh, Mauro Bisson, Jean Kossaifi, Karthik Kashinath,
736 Anima Anandkumar, William D. Collins, Michael S. Pritchard, and Alexander Keller. Four-
737 castnet 3: A geometric approach to probabilistic machine-learning weather forecasting at
738 scale, 2025. URL <https://arxiv.org/abs/2507.12144>.

739 Noah D. Brenowitz, Yair Cohen, Jaideep Pathak, Ankur Mahesh, Boris Bonev, Thorsten Kurth,
740 Dale R. Durran, Peter Harrington, and Michael S. Pritchard. A practical probabilistic
741 benchmark for ai weather models. *Geophysical Research Letters*, 52(7), April 2025. ISSN
742 1944-8007. doi: 10.1029/2024gl113656. URL <http://dx.doi.org/10.1029/2024GL113656>.

743 Lei Chen, Xiaohui Zhong, Feng Zhang, Yuan Cheng, Yinghui Xu, Yuan Qi, and Hao Li.
744 FuXi: a cascade machine learning forecasting system for 15-day global weather forecast. *npj*
745 *Climate and Atmospheric Science*, 6(1), November 2023. ISSN 2397-3722. doi: 10.1038/
746 [s41612-023-00512-1](http://dx.doi.org/10.1038/s41612-023-00512-1). URL <http://dx.doi.org/10.1038/s41612-023-00512-1>.

747 ECMWF. aifs-single-1.1 (revision 7976552), 2025a. URL [https://huggingface.co/ecmwf/](https://huggingface.co/ecmwf/aifs-single-1.1)
748 [aifs-single-1.1](https://huggingface.co/ecmwf/aifs-single-1.1).

749 ECMWF. Open data. 2025b. doi: 10.21957/OPEN-DATA. URL [https://www.ecmwf.int/](https://www.ecmwf.int/en/forecasts/datasets/open-data)
750 [en/forecasts/datasets/open-data](https://www.ecmwf.int/en/forecasts/datasets/open-data).

751 European Centre for Medium-Range Weather Forecasts. Aifs 1.1.0: Permanent archive of
752 checkpoints and source code for training and inference, 2025. URL [https://zenodo.org/](https://zenodo.org/doi/10.5281/zenodo.17349820)
753 [doi/10.5281/zenodo.17349820](https://zenodo.org/doi/10.5281/zenodo.17349820).

754 Gregory J Hakim and Sanjit Masanam. Dynamical tests of a deep-learning weather prediction
755 model. *Artificial Intelligence for the Earth Systems*, 2024.

756 Paula Harder, Alex Hernandez-Garcia, Venkatesh Ramesh, Qidong Yang, Prasanna Sattigeri,
757 Daniela Szwarcman, Campbell Watson, and David Rolnick. Hard-constrained deep learning
758 for climate downscaling, 2024. URL <https://arxiv.org/abs/2208.05424>.

759 H. Hersbach, B. Bell, P. Berrisford, et al. The ERA5 global reanalysis. *QJ R Meteorol Soc*,
760 146:1999–2049, 2020. doi: 10.1002/qj.3803.

761 George J. Huffman, Erich F. Stocker, David T. Bolvin, Eric J. Nelkin, and Jackson Tan. GPM
762 IMERG Final Precipitation L3 1 day 0.1 degree x 0.1 degree V07, 2023. URL <https://doi.org/10.5067/GPM/IMERGDF/DAY/07>. Accessed: 24/07/2025.

763

764 Ian T Jolliffe and David B Stephenson, editors. *Forecast verification*. Wiley-Blackwell, Hoboken,
765 NJ, 2 edition, December 2011.

766 R. Keisler. Forecasting global weather with graph neural networks. *arXiv preprint*
767 *arXiv:2202.07575*, Feb 15 2022.

768 Chris Kent, Adam A Scaife, Nick J Dunstone, Doug Smith, Steven C Hardiman, Tom Dunstan,
769 and Oliver Watt-Meyer. Skilful global seasonal predictions from a machine learning weather
770 model trained on reanalysis data. *Npj Clim. Atmos. Sci.*, 8(1), August 2025.

771 Remi Lam, Alvaro Sanchez-Gonzalez, Matthew Willson, Peter Wirnsberger, Meire Fortu-
772 nato, Ferran Alet, Suman Ravuri, Timo Ewalds, Zach Eaton-Rosen, Weihua Hu, Alexan-
773 der Merose, Stephan Hoyer, George Holland, Oriol Vinyals, Jacklynn Stott, Alexander
774 Pritzel, Shakir Mohamed, and Peter Battaglia. Learning skillful medium-range global
775 weather forecasting. *Science*, 382(6677):1416–1421, December 2023. ISSN 1095-9203. doi:
776 10.1126/science.adi2336. URL <http://dx.doi.org/10.1126/science.adi2336>.

777 Simon Lang, Mihai Alexe, Matthew Chantry, Jesper Dramsch, Florian Pinault, Baudouin
778 Raoult, Mariana C. A. Clare, Christian Lessig, Michael Maier-Gerber, Linus Magnusson,
779 Zied Ben Bouallègue, Ana Prieto Nemesio, Peter D. Dueben, Andrew Brown, Florian Pap-
780 penberger, and Florence Rabier. AIFS – ECMWF’s data-driven forecasting system. *arXiv*
781 *preprint arXiv:2406.01465*, 2024a. URL <https://arxiv.org/abs/2406.01465>.

782 Simon Lang, Mihai Alexe, Mariana C. A. Clare, Christopher Roberts, Rilwan Adewoyin,
783 Zied Ben Bouallègue, Matthew Chantry, Jesper Dramsch, Peter D. Dueben, Sara Hahner,
784 Pedro Maciel, Ana Prieto-Nemesio, Cathal O’Brien, Florian Pinault, Jan Polster, Baudouin
785 Raoult, Steffen Tietsche, and Martin Leutbecher. AIFS-CRPS: Ensemble forecasting using
786 a model trained with a loss function based on the continuous ranked probability score. *arXiv*
787 *preprint arXiv:2412.15832*, 2024b. URL <https://arxiv.org/abs/2412.15832>.

788 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International*
789 *Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=Bkg6RiCqY7>.

791 Paulius Micikevicius, Sharan Narang, Jonah Alben, Gregory Diamos, Erich Elsen, David Gar-
792 cia, Boris Ginsburg, Michael Houston, Oleksii Kuchaiev, Ganesh Venkatesh, and Hao Wu.
793 Mixed precision training, 2018. URL <https://arxiv.org/abs/1710.03740>.

794 J. Muñoz Sabater, E. Dutra, A. Agustí-Panareda, C. Albergel, G. Arduini, G. Balsamo,
795 S. Boussetta, M. Choulga, S. Harrigan, H. Hersbach, B. Martens, D. G. Miralles, M. Piles,
796 N. J. Rodríguez-Fernández, E. Zsoter, C. Buontempo, and J.-N. Thépaut. Era5-land: a
797 state-of-the-art global reanalysis dataset for land applications. *Earth System Science Data*,
798 13(9):4349–4383, 2021. doi: 10.5194/essd-13-4349-2021. URL <https://essd.copernicus.org/articles/13/4349/2021/>.

800 NASA Earth Observatory. Rainy, stormy days in queensland. NASA Earth Observatory,
801 Visible Earth, February 2025. URL <https://earthobservatory.nasa.gov/images/153914/rainy-stormy-days-in-queensland>. Image created by Michala Garrison using IMERG
802 data; story by Kathryn Hansen.

804 Thomas Nils Nipen, Håvard Homleid Haugen, Magnus Sikora Ingstad, Even Marius Nordhagen,
805 Aram Farhad Shafiq Salihi, Paulina Tedesco, Ivar Ambjørn Seierstad, Jørn Kristiansen,
806 Simon Lang, Mihai Alexe, Jesper Dramsch, Baudouin Raoult, Gert Mertes, and Matthew
807 Chantry. Regional data-driven weather modeling with a global stretched-grid, 2024. URL
808 <https://arxiv.org/abs/2409.02891>.

809 J. Pathak, S. Subramanian, P. Harrington, S. Raja, A. Chattopadhyay, M. Mardani, T. Kurth,
810 D. Hall, Z. Li, K. Azizzadenesheli, and P. Hassanzadeh. FourCastNet: A global data-
811 driven high-resolution weather model using adaptive fourier neural operators. *arXiv preprint*
812 *arXiv:2202.11214*, Feb 22 2022.

813 Uwe Pfeifroth, Steffen Kothe, Jaqueline Drücke, Jörg Trentmann, Marc Schröder, Nathalie
814 Selbach, and Rainer Hollmann. Surface radiation data set - heliosat (sarah) - edition 3, 2023.
815 URL https://wui.cmsaf.eu/safira/action/viewDoiDetails?acronym=SARAH_V003.

- 816 Stephan Rasp, Stephan Hoyer, Alexander Merose, Ian Langmore, Peter Battaglia, Tyler
817 Russell, Alvaro Sanchez-Gonzalez, Vivian Yang, Rob Carver, Shreya Agrawal, Matthew
818 Chantry, Zied Ben Bouallegue, Peter Dueben, Carla Bromberg, Jared Sisk, Luke Barrington,
819 Aaron Bell, and Fei Sha. Weatherbench 2: A benchmark for the next generation
820 of data-driven global weather models. *Journal of Advances in Modeling Earth Systems*,
821 16(6):e2023MS004019, 2024. doi: <https://doi.org/10.1029/2023MS004019>. URL <https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2023MS004019>. e2023MS004019
822 2023MS004019.
- 824 Mark J. Rodwell, David S. Richardson, Tim D. Hewson, and Thomas Haiden. A new equitable
825 score suitable for verifying precipitation in numerical weather prediction. *Quarterly Journal*
826 *of the Royal Meteorological Society*, 136(650):1344–1363, 2010. doi: <https://doi.org/10.1002/qj.656>. URL <https://rmets.onlinelibrary.wiley.com/doi/abs/10.1002/qj.656>.
- 828 John S. Schreck, Yingkai Sha, William Chapman, Dhamma Kimpara, Judith Berner, Seth
829 McGinnis, Arnold Kazadi, Negin Sobhani, Ben Kirk, Charlie Becker, Gabrielle Gantos,
830 and David John Gagne II. Community research earth digital intelligence twin: a scalable
831 framework for ai-driven earth system modeling. *npj Climate and Atmospheric Science*, 8(1),
832 June 2025. ISSN 2397-3722. doi: 10.1038/s41612-025-01125-6. URL <http://dx.doi.org/10.1038/s41612-025-01125-6>.
- 834 Yingkai Sha, John S. Schreck, William Chapman, and David John Gagne II. Investi-
835 gating the use of terrain-following coordinates in ai-driven precipitation forecasts. *Geo-*
836 *physical Research Letters*, 52(20):e2025GL118478, 2025a. doi: <https://doi.org/10.1029/2025GL118478>. URL <https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2025GL118478>. e2025GL118478 2025GL118478.
- 839 Yingkai Sha, John S. Schreck, William Chapman, and David John Gagne II. Improving ai
840 weather prediction models using global mass and energy conservation schemes. *Journal of*
841 *Advances in Modeling Earth Systems*, 17(11):e2025MS005138, 2025b. doi: <https://doi.org/10.1029/2025MS005138>. URL <https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2025MS005138>. e2025MS005138 2025MS005138.
- 844 Akshay Subramaniam, Dale Durran, David Pruit, Nathaniel Cresswell-Clay, and William Yik.
845 Imposing the fundamental dynamical constraint of hydrostatic balance to improve global ml
846 weather prediction, 2025. URL <https://arxiv.org/abs/2506.08285>.
- 847 Haley Thiem and Nicole Collins. Historic January 2025 snowstorm in the south-
848 ern US, 2025. URL <https://www.climate.gov/news-features/event-tracker/historic-january-2025-snowstorm-southern-us>.
- 850 N. P. Wedi. Increasing the horizontal resolution in numerical weather prediction and climate
851 simulations: illusion or panacea? *Philosophical Transactions of the Royal Society A*, 372,
852 2014. doi: 10.1098/rsta.2013.0289.
- 853 Jasper S. Wijnands, Michiel Van Ginderachter, Bastien François, Sophie Buurman, Piet Term-
854 onia, and Dieter Van den Bleeken. A comparison of stretched-grid and limited-area modelling
855 for data-driven regional weather forecasting, 2025. URL <https://arxiv.org/abs/2507.18378>.
- 857 Daniel S Wilks. *Statistical methods in the atmospheric sciences*. Elsevier Science Publishing,
858 Philadelphia, PA, 4 edition, June 2019.

859 Hao Zuo, Magdalena Alonso-Balmaseda, Eric de Boisseson, Philip Browne, Marcin Chrust,
860 Sarah Keeley, Kristian Mogensen, Charles Pelletier, Patricia de Rosnay, and Toshinari
861 Takakura. Ecmwf's next ensemble reanalysis system for ocean and sea ice: Oras6. *ECMWF*
862 *Newsletter*, (180):30–36, 07/2024 2024. doi: 10.21957/hzd5y821lk.