

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

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

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

10 ²European Centre for Medium-Range Weather Forecasts, Robert-Schuman-Platz 3,
11 53175 Bonn, Germany

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

13 September 2025

14 **Abstract**

15 We present ~~an update to version 1.1.0 of ECMWF’s machine-learned weather forecasting~~
16 ~~model AIFS Single with several key improvements. The model now incorporates physical~~
17 ~~consistency constraints through bounding layers, an updated training schedule, and an expanded~~
18 ~~set of variables. The physical constraints substantially improve precipitation forecasts and~~
19 ~~the new variables show a high level of Artificial Intelligence Forecasting System (AIFS Single),~~
20 ~~operational since 25 February 2025. The revised system introduces a bounding-layer framework~~
21 ~~that enforces physical constraints, such as non-negativity and internal consistency within~~
22 ~~precipitation and cloud cover variables, alongside expanded training data, revised loss weighting,~~
23 ~~and an extended set of surface and atmospheric variables. Overall skill improves by 4–6% in~~
24 ~~the upper air and near-surface variables without degradation of spatial variability. A controlled~~
25 ~~comparison shows that training data expansion is the dominant source of upper-air skill gains,~~
26 ~~highlighting the importance of frequent model updates. The bounding framework delivers~~
27 ~~the largest precipitation improvements, up to 12% and an approximately one-day advantage~~
28 ~~using a categorical measure of skill. Upper-air headline scores also show improvement over~~
29 ~~the previous AIFS version. The AIFS has been fully operational at ECMWF since the 25th of~~
30 ~~February 2025. We further show that enforcing precipitation non-negativity resolves a gradient~~
31 ~~ambiguity at the zero-precipitation boundary under MSE training, explaining the reduction in~~
32 ~~drizzle bias and the improvements in precipitation.~~

*equal contribution

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

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

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

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

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

In this paper we begin by outlining the training setup of the model and how this differs from

82 the previous AIFS version. Then we motivate and describe the new bounding strategy to make
83 the model forecast more physically consistent. We demonstrate the improved performance of
84 the revised AIFS version via evaluation results and selected case studies. We conclude by
85 summarizing main results and future work in the discussion and conclusions.

86 2 Training

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

90 The model operates on a reduced Gaussian grid, (N320, approximately 0.25° resolution).
91 The processor (or hidden) grid is an O96 octahedral reduced Gaussian grid (Wedi (2014)) with
92 40,320 grid points, approximately 1° resolution, and consists of 16 processor layers.

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

97 2.1 Training Schedule

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

107 Pre-training is performed on ERA5 data covering the years 1979–2022 (compared to 1979–2020
108 in the previous AIFS version), using a cosine learning rate (LR) schedule, a batch size of 16,
109 and a total of 260,000 training steps. The LR is linearly increased from 0 to 5×10^{-4} during
110 the first 1,000 steps, then annealed to a minimum of 3×10^{-7} . This is followed by rollout fine-
111 tuning on ECMWF operational analysis from 2016 to 2022, also using a cosine LR schedule
112 and batch size of 16, for approximately 7,900 steps (equivalent to one epoch per rollout step).
113 The LR started at 1.28×10^{-5} and is annealed to the same minimum value of 3×10^{-7} . The
114 rollout length is initially set to 6 hours (1 step) and progressively increased by one step per
115 epoch up to 72 hours (12 steps), following the approach of Lam et al. (2023) and Lang et al.
116 (2024a). We used the AdamW optimizer (Loshchilov and Hutter, 2019) with β coefficients of
117 0.9 and 0.95. Here, the rollout dataset is extended to eight years of operational IFS analysis
118 (2016–2022), compared with only two years (2019–2020) in the previous AIFS version.

119 2.2 Variables used in training

120 The variables used in the new AIFS version are listed in Table 1. As in AIFS 0.2.1, the upper
121 atmosphere is represented by geopotential, horizontal wind components, specific humidity,
122 and temperature at 13 pressure levels: 50, 100, 150, 200, 250, 300, 400, 500, 600, 700, 850,
123 925, and 1000 hPa. Newly introduced variables are marked with *. We have increased the
124 characterization of the land surface in the model by including new prognostic variables of soil
125 moisture at levels 1 and 2 (swvl1 and swvl2), and soil temperature at levels 1 and 2 (stl1

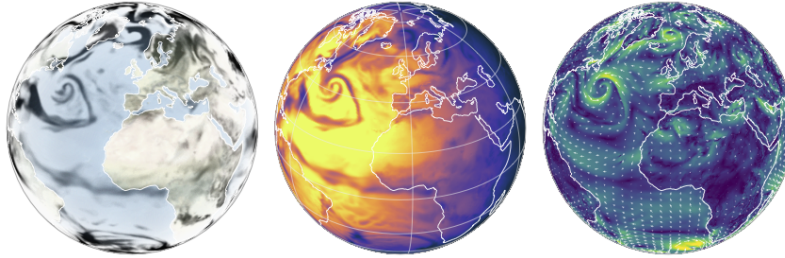


Figure 1: A selection of new variables available from the revised AIFS Single forecasts: cloud cover (left), surface solar radiation (centre), and 100 m wind speed/direction (right). The consistency between these new variables is clear, with areas of higher cloud cover corresponding to lower solar radiation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

126 and stl2), important for drought monitoring and forecasting. A notion of hydrology has been
 127 included with runoff (ro), forecast as a diagnostic variable. A second set of variables, related
 128 to energy forecasting and clouds, adds real value to the model’s utility. These are forecast
 129 diagnostically and include the 100-metre wind components (100u and 100v), surface solar and
 130 thermal radiation (ssrd and strd), and cloud cover at various levels (tcc, hcc, mcc, lcc). Finally,
 131 snowfall (sf) has been added to complement the set of total precipitation-related variables. An
 132 illustration of a selection of these variables can be seen in the forecast presented in Figure 1,
 133 where the consistency between these new variables is clear, with areas of higher cloud cover
 134 corresponding to lower solar radiation at the surface and consistent weather patterns for 100-
 135 metre winds. These new variables are sourced from the ERA5 reanalysis and IFS operational
 136 data archive, in line with those used in the previous AIFS version (0.2.1).

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

154 AIFS is trained using data parallelism with a batch size of 16, while each model instance is
 155 distributed across four GPUs within a single node (Lang et al., 2024a). Training was conducted
 156 on the European supercomputer Leonardo (EuroHPC), hosted and managed by Cineca, on
 157 64GB A100 GPUs. Mixed-precision training is used (Micikevicius et al. (2018)), and the full
 158 process takes approximately three days. A 10-day forecast can be produced in about 2 minutes
 159 and 30 seconds on a single A100 40GB GPU, including data input and output.

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.

3 Enforcing Model Constraints

Machine-learned forecast models for numerical weather prediction show very good forecast skill, yet they are prone to producing outputs that violate known physical laws or expected statistical consistency. Unlike traditional numerical models, which are governed by equations ensuring mass conservation, positivity, or energy bounds, machine-learned forecast models lack such guarantees by default. As a result, physically implausible outputs, such as negative precipitation, can emerge. We show that incorporating constraints into the model design to enforce physical realism improves forecast skill. In this section, we first identify specific issues in the output of the previous AIFS version related to total precipitation, and then introduce a simple yet effective method to bound the model outputs using activation functions. The proposed method is not restricted to total precipitation but can be equally applied to other variables.

3.1 Lack of Physical Realism in Precipitation Forecasts

The previous AIFS version suffers from significant drawbacks in forecasting precipitation. Most notably, the model's output is not constrained, leading to a frequent occurrence of negative values. This is illustrated in Figure 2, which compares the 24-hour accumulated total precipitation forecasts from the previous AIFS version, the revised version, ~~and an estimate derived from the short-range GraphCast (Lam et al., 2023), FuXi (Chen et al., 2023) and an IFS (47r3) 6-hour forecasts~~24-hour forecast, for the run initialized on 01/06/2023 at 00:00 UTC and valid at 02/06/2023 00:00 UTC. The previous AIFS ~~shows model and GraphCast show~~ spurious negative precipitation values ~~and an excess of light rainfall~~, which are largely corrected in the revised AIFS. While negative values can be clipped to zero at inference time (as is done in FuXi in this figure and thus non visible), their presence highlights a lack of physical consistency in the model. ~~This issue is also present in other machine-learned weather forecast models, such as GraphCast (Lam et al., 2023), which is similarly unconstrained.~~

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

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

Convective precipitation forecasts also exhibit similar shortcomings. In addition, there is a further lack of physical consistency. Convective precipitation represents the part of the total precipitation that originates from convection, and therefore should always be less than or equal to the total. Figure 4 shows the previous AIFS 24-hour accumulated forecasts of total

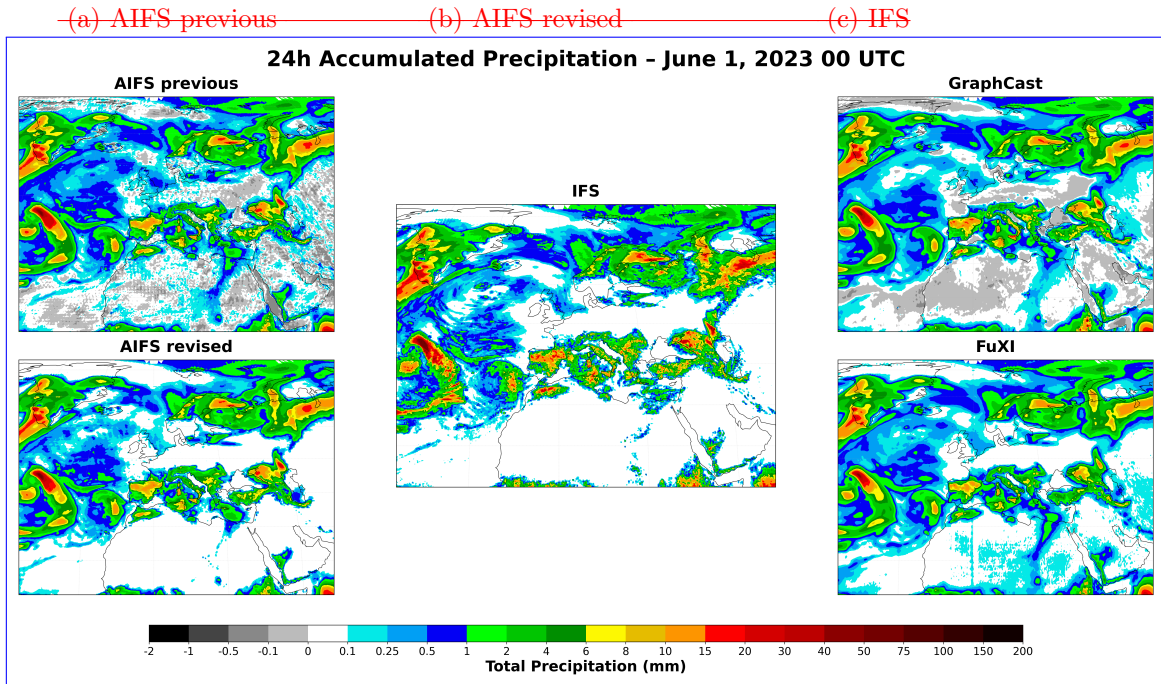


Figure 2: Comparison of 24-hour total precipitation accumulation from ~~the previous AIFS, the revised AIFS and an estimate derived from the short-range IFS (47r3) 6-hour forecasts, five forecasting systems~~ for the forecast issued at 01/06/2023 00:00 UTC and valid at 02/06/2023 00:00 UTC: previous AIFS, revised AIFS, operational IFS, GraphCast, and FuXi. The previous AIFS shows spurious negative precipitation values, GraphCast, and FuXi all exhibit an excess of light rainfall, which are largely corrected in the revised AIFS characteristic biases of ML weather models. The revised version therefore 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.

207 and convective precipitation for 02/06/2023. The map displaying the difference between the
 208 two reveals frequent cases in which convective precipitation exceeds total precipitation, which
 209 should not occur.

210 The CREDIT platform Schreck et al. (2025) has recently been used to explore physically
 211 informed constraints for addressing drizzle bias; Sha et al. (2025b) implemented global mass
 212 and energy conservation schemes as modular constraints within FuXi and demonstrated a
 213 direct reduction of drizzle bias; a companion study Sha et al. (2025a) further showed that
 214 incorporating terrain-following (hybrid sigma-pressure) can improve extreme precipitation forecasts.

215
 216 Here, we address the drizzle and negative precipitation issue through simplified intervention:
 217 enforcing only the physically admissible output range via a hard-constraint. This approach is
 218 described in Section 3.2. In Section 4.1, we show that this minimal architectural modification
 219 fundamentally reshapes the loss landscape in the vicinity of zero precipitation, eliminating
 220 gradient ambiguity and substantially reducing light-precipitation bias.

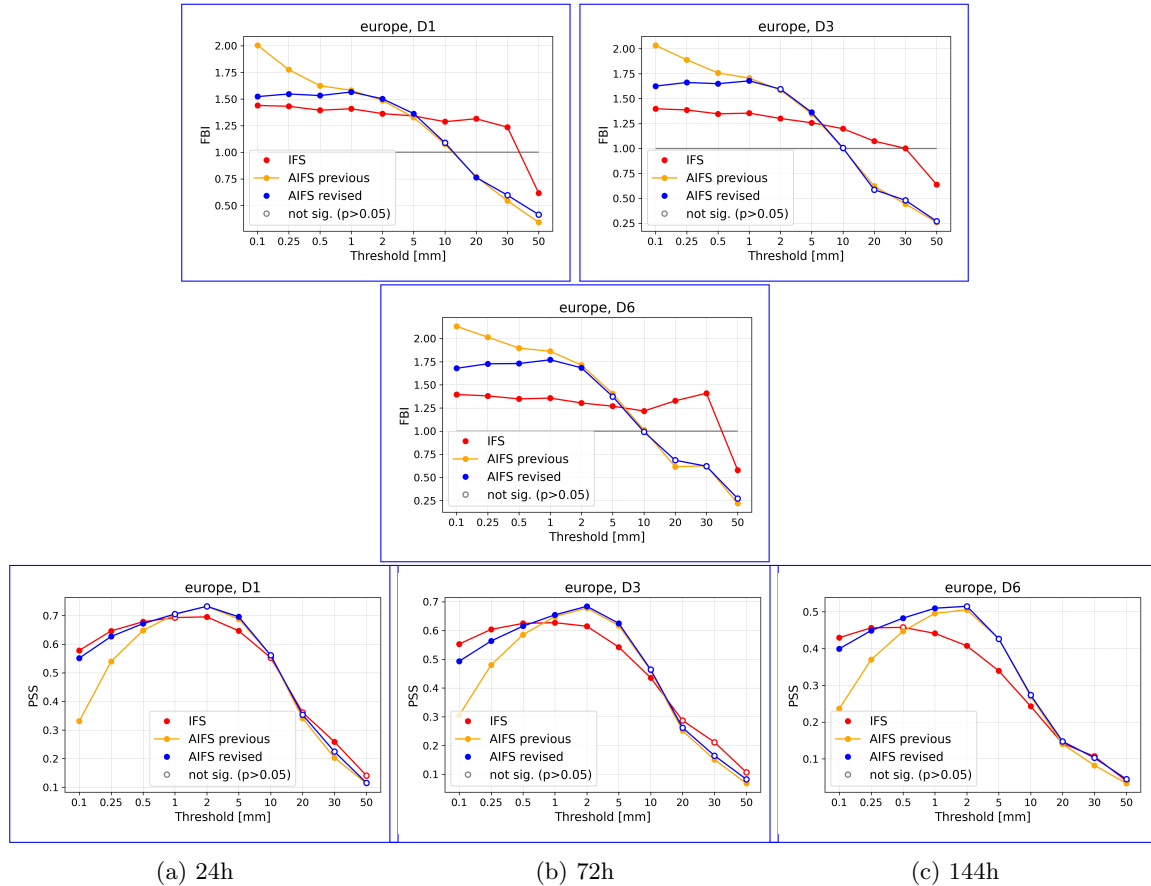


Figure 3: Comparison of IFS-Frequency Bias Index (red) FBI, revised AIFS (blue) FBI, and previous AIFS-Peirce Skill Score (orange) PSS, bottom) for 24-hour accumulated precipitation over Europe as a function of threshold, at forecast steps 24h, 72h, and 144h for 2023. Top row: Frequency Bias Index (FBI) left to right; Bottom row: Peirce Skill Score. Scores are averaged over all initialisation dates in 2023. Filled markers indicate that the difference relative to the previous AIFS version is statistically significant (PSS 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 the AIFS predicts light precipitation in excess.

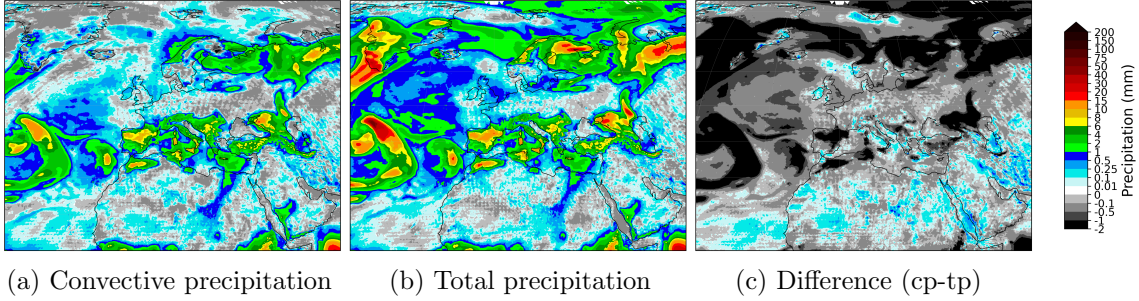


Figure 4: Comparison of 24-hour total and convective precipitation forecast from the previous AIFS version, together with a map showing the difference between the two of them for the forecast issued at 01/06/2023 00:00 UTC and valid at 02/06/2023 00:00 UTC. Positive values (coloured regions) in the difference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

3.2 Bounding the Outputs with Activation Functions

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

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

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

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

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

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

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

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

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

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

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

Bounding Type	Range	Variables
ReluBounding	$[0, \infty)$	$tp, ro, tcw, ssrd, q(50-1000 \text{ hPa})$
HardtanhBounding	$[0, 1]$	$tcc, swvl1, swvl2$
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.

275 4 Evaluation

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

279 The revised AIFS version delivers highly skilled forecasts, as shown by anomaly corre-
 280 lation scores for 2023 in the Northern Hemisphere (Figure 8). In the medium range (3-10

281 days), AIFS outperforms the IFS by 12 to 24 hours in skill. Forecast skill is also clearly
 282 improved compared to the previous AIFS version. This performance gain can be attributed
 283 to ~~more training data~~ and the combined effect of increased training data, improvements in
 284 rollout fine-tuning. ~~Here, we,~~ the implementation of output bounding, and the inclusion of
 285 new prognostic variables. To quantify the specific contribution of expanded training data, we
 286 present a controlled comparison in Figure 5. We verify against the operational IFS analysis,
 287 which is also used to initialise the forecasts.

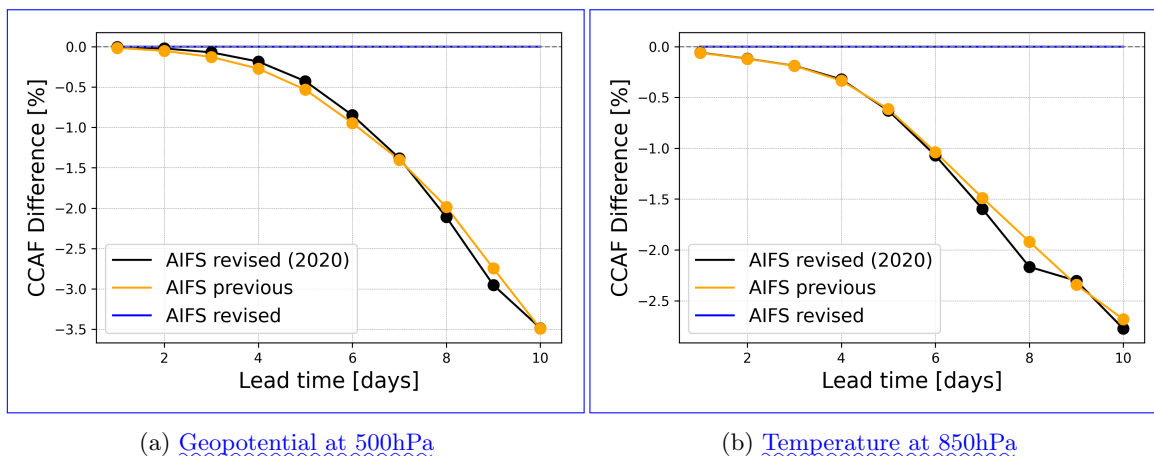


Figure 5: Anomaly correlation skill score difference for Geopotential at 500hPa and Temperature at 850hPa for 2023. This controlled comparison shows: (1) AIFS revised model (full system with all modifications), (2) AIFS revised trained with limited data (ERA5 up to 2020, rollout fine-tuning 2019-2020 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, $p < 0.05$).

288 As shown in Figure 5, the expanded training dataset contributes to the most important
 289 portion of the overall performance gain. This indicates that data availability (ERA5 extended
 290 to 2022 and rollout fine-tuning expanded from 2019-2020 to 2016-2022) plays a major role. The
 291 remaining improvement stems from other system modifications, including rollout fine-tuning
 292 schedules, output bounding layers, and expanded prognostic variables. Due to the high computational
 293 cost, a detailed ablation study to isolate the impact of each individual modification beyond data
 294 expansion was not performed; thus, the observed improvements represent the cumulative result
 295 of these integrated system updates. It should be noted that the close agreement between AIFS
 296 revised (2020 data) and AIFS previous in ACC should be interpreted with caution, as these
 297 configurations differ in their training protocols: AIFS previous includes a rollout fine-tuning
 298 phase on ERA5 which AIFS revised (2020) does not, and uses only 2 years (2019-2020) of
 299 operational data for final rollout fine-tuning compared to 6 years (2016-2022) in the full revised
 300 version. Furthermore, similar ACC scores do not imply equivalent forecast quality. As shown
 301 in Figure 6, AIFS revised (2020) exhibits less mesoscale smoothing than AIFS previous despite
 302 comparable ACC, indicating that the changes introduced in the revised system do contribute
 303 positively to forecast quality in ways not fully captured by ACC alone.

304 Additionally, imposing a minimum on the loss weights in the stratosphere leads to significant
305 improvements in the data-driven forecasts at 100 and 50 hPa (Figure 9). For temperature at
306 100hPa, the new version of the AIFS outperforms the IFS, while for 50hPa wind speed, the gap
307 in skill between the previous version of AIFS and the IFS in the stratosphere is significantly
308 reduced.

309 Forecast skill for key surface variables, such as 2-metre temperature and 10-metre wind
310 speed, verified against SYNOP observations, is similarly improved (Figure 10). Overall, the
311 new AIFS version exhibits improvements of around 4–6 % across all variables, lead times, and
312 pressure levels relative to the previous AIFS version, as shown in the scorecard presented in
313 Figure 7. The performance of the model for tropical cyclone prediction is similar to that of the
314 previous version (see Lang et al. (2024a)), with some small improvements to track position.
315 ~~As a design choice, rollout fine-tuning was configured to ensure that the field smoothness~~
316 ~~characteristics remain consistent with those of the~~ The training configuration, including a
317 maximum rollout length of 12 (72 hours), was retained from the previous AIFS version. ~~This~~
318 ~~was confirmed by spectral analysis (not shown here)~~, as shown in Section 2.1. This parameter
319 is known to influence spectral characteristics, with longer rollouts leading to enhanced damping.
320 No explicit tuning was performed to target spectral behaviour.

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

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

338 The forecasting skill of the model with respect to 24-hour accumulated total precipitation
339 is significantly improved. The new AIFS version is compared against both the previous AIFS
340 version and the operational IFS (cycles 47r3 and 48r1) in Figure 12. The Stable Equitable
341 Error in Probability Space (SEEPS) skill score (Rodwell et al. (2010)) is used as the primary
342 verification metric, with 24-hour accumulated precipitation SYNOP observations serving as the
343 reference. Results show a consistent and statistically significant improvement across all lead
344 times and in the ~~three main global regions: the Northern Hemisphere, Northern Hemisphere~~
345 ~~and the Southern Hemisphere, and the tropics.~~ The revised AIFS demonstrates approximately
346 a one-day gain in forecast skill relative to both IFS and the previous AIFS version. The forecast
347 fields also exhibit noticeable improvements, as illustrated in Figure 2. The new version of the
348 AIFS produces no negative values in the output and substantially reduces light precipitation,
349 aligning more closely with the 24-hour total precipitation accumulation fields derived from the
350 IFS operational short-range forecasts.

351 Figure 3 reveals where the improvement originates. The Frequency Bias Index (FBI)
352 ~~and~~, Wilks 2019), defined as the ratio of predicted to observed event frequency at a given
353 threshold ($FBI = (H + FA)/(H + M)$, where H are hits, FA false alarms, and M misses),

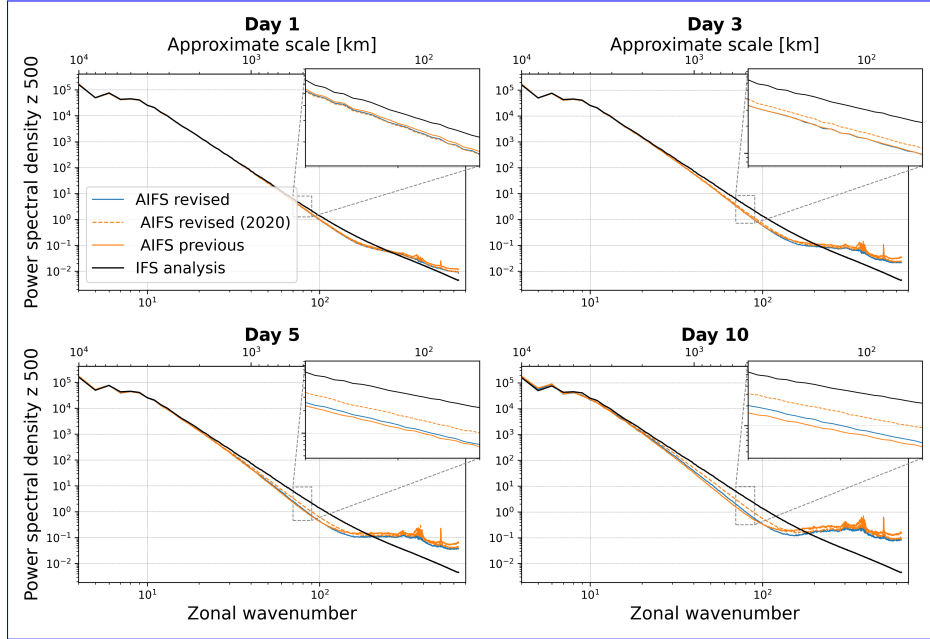


Figure 6: [Z500 power spectral density as a function of zonal wavenumber \(bottom axis\) and approximate horizontal scale in km \(top axis\) for forecast lead times Day 1, 3, 5, and 10 during JJA 2023. Spectra from the revised AIFS \(blue\), AIFS revised trained with limited data \(ERA5 up to 2020, rollout fine-tuning 2019-2020 only\) 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.](#)

354 [and the Peirce Skill Score \(PSS\)–are–, also known as the Hanssen–Kuipers discriminant;](#)
 355 [Jolliffe and Stephenson 2011\), defined as the difference between the probability of detection](#)
 356 [and the probability of false detection \(\$PSS = H/\(H + M\) - FA/\(FA + CN\)\$, where \$CN\$ are](#)
 357 [correct negatives\), are shown for the Northern Hemisphere for different thresholds. The pre-](#)
 358 [vious AIFS version exhibits a strong tendency to over-predict light precipitation events \(< 1](#)
 359 [mm\) across all lead times, as shown by the FBI. This bias is substantially corrected due to the](#)
 360 [bounding \(see Section 4.1\) in the revised AIFS.](#)

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

367 This under-prediction plays an important role in the metrics concerning more extreme
 368 events, since both the previous and the revised AIFS models underperform IFS for thresholds
 369 exceeding 10mm in terms of PSS, but remains competitive. This suggests that although the
 370 AI models predict fewer high-intensity events, their predictions are more accurate when they

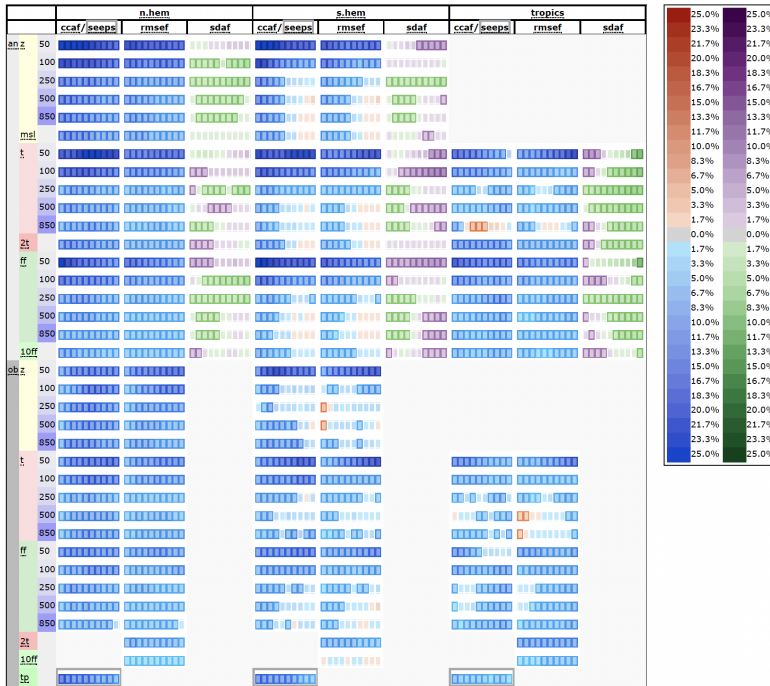


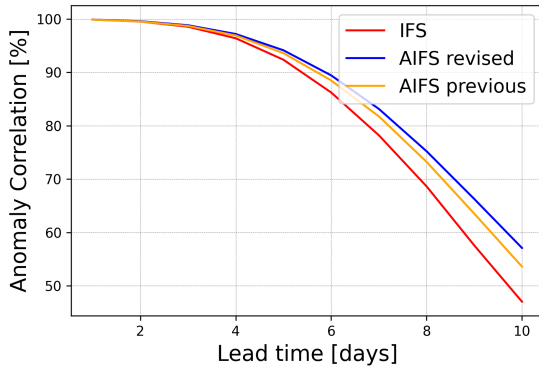
Figure 7: Scorecard comparing forecast scores of AIFS revised versus the previous AIFS version for the whole year of 2023. Forecasts are initialised on 00 and 12 UTC. Relative score changes are shown as function of lead time (day 1 to 10) for northern extra-tropics (n.hem), southern extra-tropics (s.hem) 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).

371 do occur. Finally, the revised AIFS shows a marginal improvement in terms of PSS compared
 372 against the previous AIFS version, possibly due to improvements in the learning-rate scheduling
 373 used for fine-tuning and additional training data.

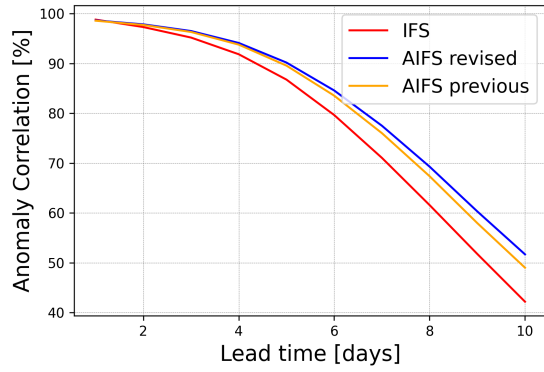
374 4.1 Evaluating the effects of bounding on total precipitation

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

381 Other factors that might improve the precipitation forecast skill in the revised AIFS version

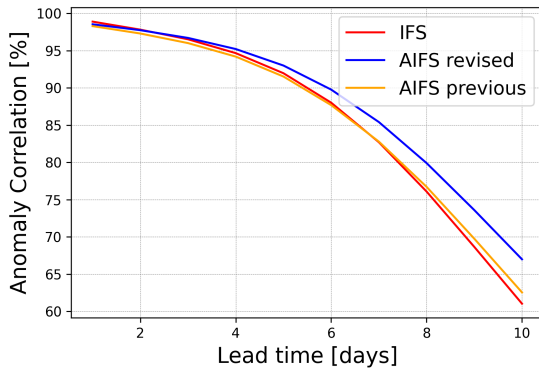


(a) Geopotential at 500hPa

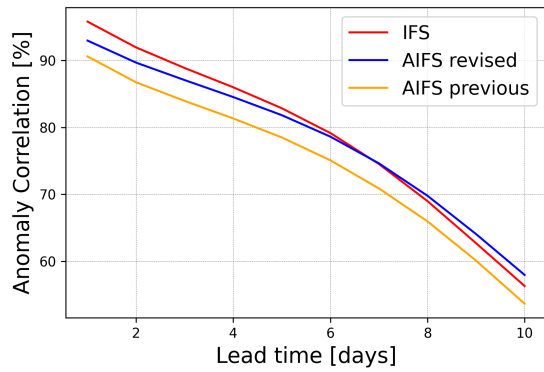


(b) Temperature at 850hPa

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.

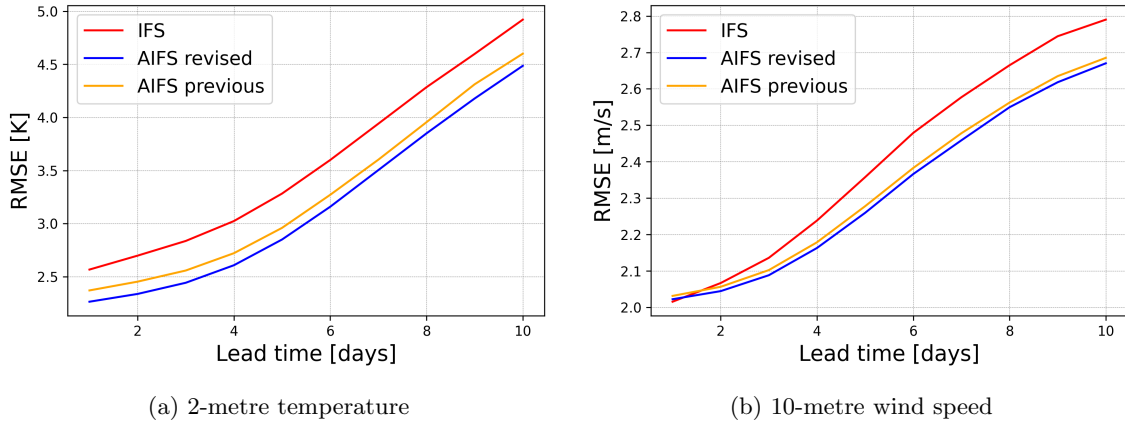


Figure 10: RMSE scores for 2-metre temperature and 10-metre wind speed computed against SYNOP observations over the Northern Hemisphere. The revised AIFS version shows improvement when compared to the previous version of the AIFS.

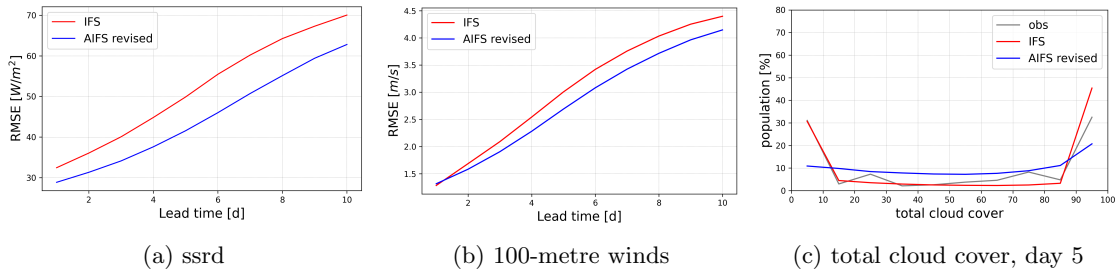


Figure 11: Forecast RMSE computed against operational IFS analysis and distribution comparison for new variables. (a) Surface solar radiation downwards RMSE for March–May (MAM) 2023, (b) 100-metre wind speed RMSE for the full year 2023, (c) Total cloud cover distribution for June–August (JJA) 2023. Blue 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.

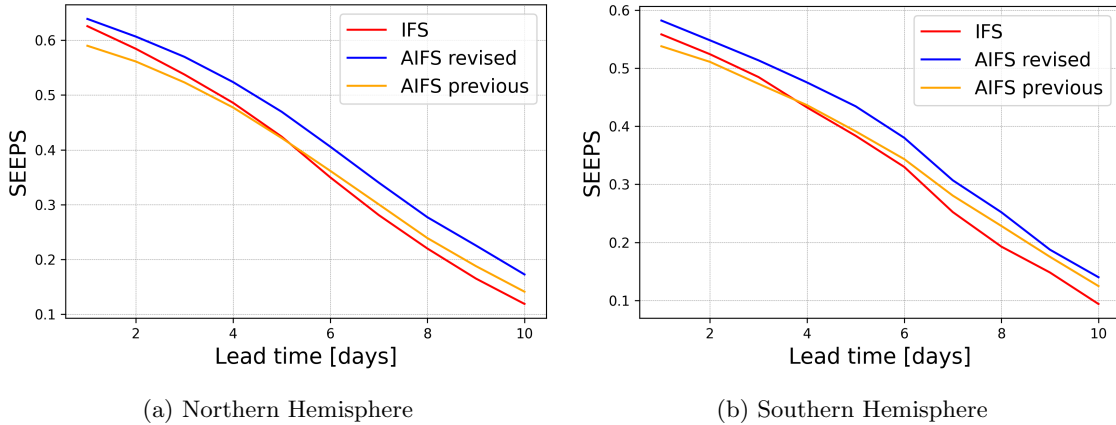


Figure 12: SEEPS skill scores for 2023 based on 24-hour accumulated precipitation from SYNOP observations, comparing the revised AIFS (blue), the previous AIFS version (orange), and the IFS (red) across different regions. Results show a consistent and statistically significant improvement across all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

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

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

399 To better understand the mechanisms ~~behind total precipitation forecasting governing total~~
 400 precipitation forecasts in the revised AIFS ~~version configuration~~, we examine the model²'s be-
 401 haviour in the negative ~~forecast space~~, revealed pre-activation space obtained by removing
 402 the final ReLU layer ~~(Figure 15). Bounding at inference.~~ Figure 15 reveals that this nominally
 403 hidden negative space is neither random nor noisy, but highly structured.

404 At first glance, bounding an output variable via ReLU has some drawbacks with a ReLU
 405 activation may appear to introduce a drawback: the negative ~~space is unconstrained since any~~
 406 ~~changes in model behaviour in the negative space are mapped~~ pre-activation space is not directly
 407 penalized, since all negative values are projected to zero before the loss is computed, which
 408 ~~means that these points evaluated.~~ In principle, changes within this region do not influence the

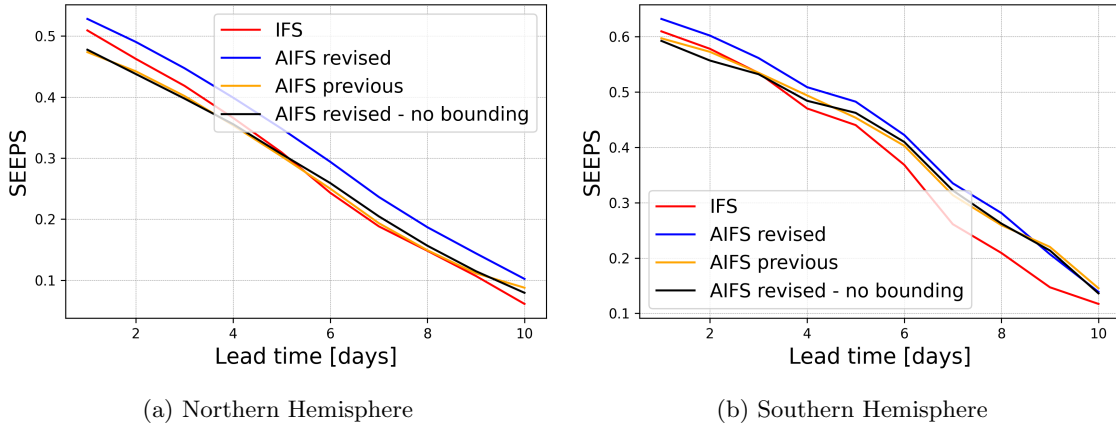


Figure 13: SEEPS skill scores for 2023 JJA comparing revised AIFS (blue), revised AIFS without bounding (black), previous AIFS (orange), and IFS (red) across different regions. The improvement observed in total precipitation forecast skill in the revised AIFS version can mainly be attributed to bounding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

409 weight ~~update. Interestingly, this hidden negative space shows updates. One might therefore~~
 410 ~~expect the negative space to be uninformative or unstable.~~
 411 ~~Instead, we observe a coherent and structured pattern. Very physically meaningful organization.~~
 412 ~~Persistently dry regions, such as the Sahara Desert, exhibit strongly negative values pre-activations,~~
 413 ~~while areas near precipitation events gradually approach zero in a smooth and continuous~~
 414 ~~manner. This suggests that the model has implicitly learned to use approaching precipitation~~
 415 ~~events transition smoothly toward zero. The model has therefore learned to encode dryness in~~
 416 ~~the negative space as a proxy for “, effectively using it as a latent representation of the “no-rain”~~
 417 ~~classification” regime.~~
 418 ~~(a) AIFS previous (b) AIFS rev. (neg. This observation motivates two fundamental questions:~~
 419 ~~(i) why does the negative pre-activation space contain coherent and physically meaningful~~
 420 ~~structure, and (ii) why does enforcing a non-negativity constraint during training improve~~
 421 ~~light-precipitation skill? We argue that the first arises from the shared latent representation of~~
 422 ~~the atmospheric state learned by the network, while the second is governed by the symmetry~~
 423 ~~properties of the MSE gradient near the zero-precipitation boundary.~~

4.1.1 Representation of Dry States in the Negative Space

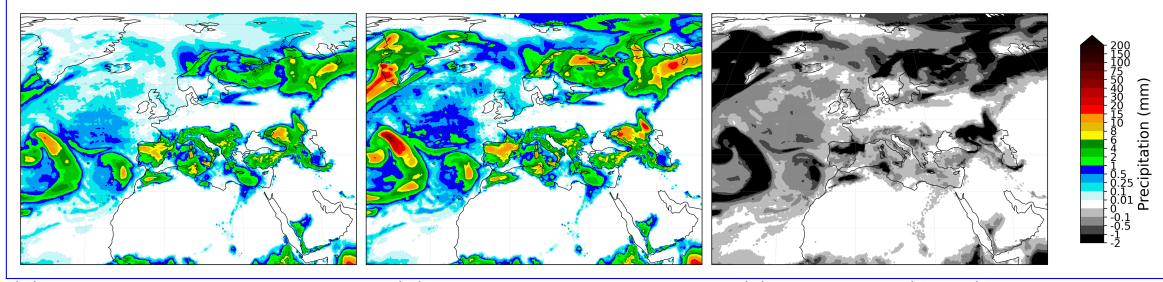
425 In this study we argue that the structure present in the negative space is an emergent feature
 426 arising from the shared representation of atmospheric states.

427 The model encodes input prognostic (\mathbf{X}_t) and forcing variables (\mathbf{F}_t) into a high-dimensional
 428 latent space (z_t) via an encoder:

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

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

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



(a) Convective precipitation (b) Total precipitation (c) Difference (cp-tp)

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.

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

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

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

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

438 Because Decoder_{tp} maps from a latent space optimized for smooth gradients (z_{t+6}), η
 439 inherits this spatial structure. The precipitation decoder head learns a smooth mapping
 440 from the latent space encoding the moisture state of the system to physical precipitation
 441 in the positive regime ($\eta > 0$), where gradients are active. Because neural networks are
 442 continuous functions biased toward smoothness, this "moisture-to-precipitation" logic naturally
 443 extrapolates into the negative regime. As moisture variables decrease, the decoder continues
 444 to output decreasing values, pushing η into the negative space.

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

450 Ultimately, this reveals that the optimization of the shared latent space is driven by the
 451 collective constraints of all output variables. In this framework, the negative pre-activation
 452 space for precipitation serves as a "saturation deficit" proxy that is kept physically consistent
 453 by the gradients flowing from prognostic moisture fields. The shared representation of the
 454 atmosphere in the latent space allows the model to maintain a sophisticated, structured representation
 455 of dryness even in the absence of direct precipitation gradients.

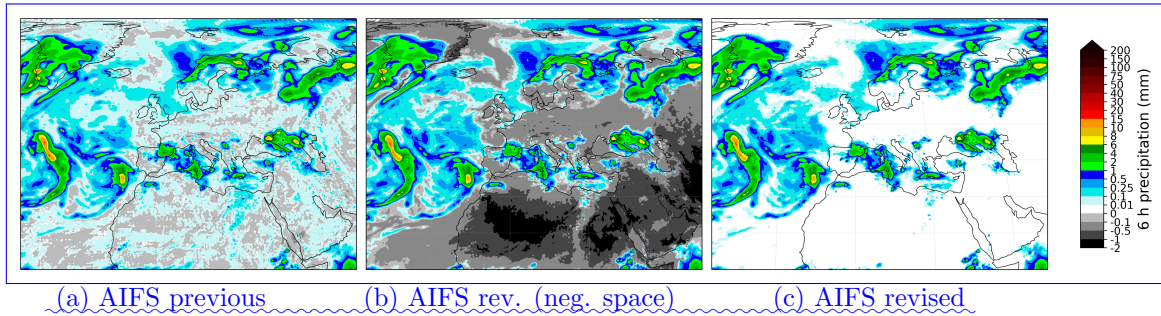


Figure 15: Comparison of 6-hour total precipitation from previous AIFS, revised AIFS without the final ReLU layer to show the negative space, and the standard revised AIFS with the final ReLU layer. Forecasts are initialised at 01/06/2023 00:00 UTC and valid at 01/06/2023 06:00 UTC. Removing the final 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.

456 To provide empirical weight to this mechanistic theory, we investigate the information
 457 content within the pre-activation space) (c)AIFS revised Comparison of 6-hour total
 458 precipitation from previous AIFS, revised AIFS without the final ReLU layer to show the
 459 negative space, and the standard revised AIFS with the final ReLU layer. Forecasts are
 460 initialised at 01/06/2023 00:00 UTC and valid at 01/06/2023 06:00 UTC. Removing the final
 461 bounding layer from the AIFS revised model reveals the behaviour of the negative space for
 462 the total precipitation variable. The model has implicitly learned to use the negative space
 463 as a proxy for “no-rain” classification. η by partitioning the model output into three distinct
 464 physical regimes: the negative (non-precipitating) space, the light precipitation regime (0–0.5
 465 mm/6h), and the moderate precipitation regime (0.5–10 mm/6h).

466 The physical consistency of convective precipitation forecast in respect to total precipitation
 467 can also be evaluated for a given forecast to assess the utility of the FractionBounding strategy
 468 used. Figure 14 presents the 24-hour total and convective precipitation accumulation together
 469 with a map showing the difference between the two for a We hypothesize that the pre-activation
 470 space η undergoes a fundamental physical decoupling as it transitions from dry to wet conditions.
 471 In the negative (non-precipitating) regime, the absence of precipitation is a deterministic
 472 function of low humidity; thus, the decoder should preserve a strong linear mapping from
 473 the prognostic moisture fields.

474 Conversely, we expect this linear correlation to weaken in the light precipitation regime
 475 ($0 < \eta < 0.5$ mm). While moisture remains a necessary condition for rain, the exact accumulation
 476 at these low intensities becomes increasingly stochastic, influenced by non-linear factors such
 477 as sub-grid scale turbulence, cloud-base evaporation, and microphysical uncertainties. These
 478 processes act as “interference,” decoupling the surface precipitation from the column moisture
 479 signal.

480 We performed a global correlation analysis on a single forecast issued at 01/06/2023 00:00
 481 UTC and valid at 02/06/2023 00:00 UTC. Unlike the previous AIFS version (Figure 4), the
 482 convective precipitation forecast is now consistent with the predicted total precipitation accumulation.
 483 . For this experiment, we utilize the AIFS revised model without the final bounding layer on
 484 tp during inference, but activated during training. We focus our analysis on the first 120 hours
 485 (5 days) of the forecast.

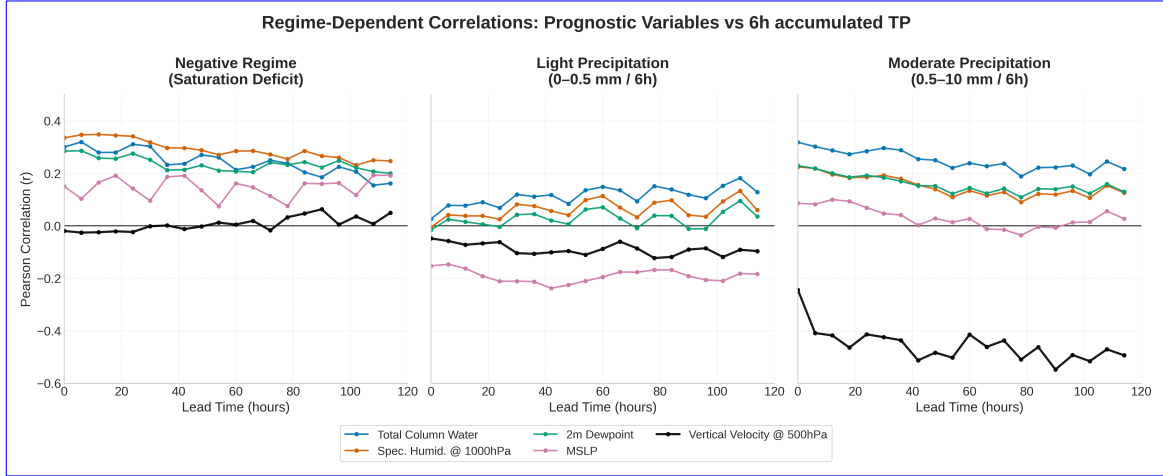


Figure 16: Regime-dependent correlations of pre-activation η (AIFS Revised), for a forecast issued the June 1, 2023 at 00 UTC. Pearson r between η and physical drivers across three regimes: (Left) Negative space ($\eta < 0$): high correlation with moisture variables (q_{1000} , TCW) identifies η as a 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.

486 We computed the Pearson correlation coefficient (r) between the pre-activation field η and
 487 five key physical drivers: Total Column Water (TCW), Specific Humidity (q_{1000}), 2m Dewpoint
 488 ($2d$), Mean Sea Level Pressure (MSLP), and mid-tropospheric Vertical Velocity (w_{500}). As
 489 shown in Figure 16, the results reveal a clear regime-dependent physical logic:

- 490 • **Negative Regime ($\eta < 0$):** We observe stable correlations ($r \approx 0.3$) with moisture
 491 variables (q_{1000} , TCW, and $2d$). This confirms that the negative space encodes a structured
 492 representation of the *saturation deficit*, kept physically consistent by gradients flowing
 493 from the prognostic moisture fields.
- 494 • **Light Precipitation ($0 < \eta \leq 0.5$ mm):** Correlation with specific humidity, 2m dewpoint
 495 and total column water is substantially reduced in this regime. The weaker relationships
 496 are consistent with a lower signal-to-noise ratio and increased sensitivity to small-scale or
 497 non-linear processes.
- 498 • **Moderate Precipitation ($1 < \eta \leq 10$ mm):** The model transitions to dynamic control.
 499 While moisture correlations remain moderate, Vertical Velocity (w_{500}) emerges as the
 500 primary physical driver ($r \approx -0.5$), illustrating the model’s reliance on large-scale ascent
 501 to produce deterministic rainfall.

502 While presented as a targeted demonstration of internal model behaviour, the consistency
 503 of these signals across lead times suggests that this regime-specific transition is a fundamental
 504 structural property of the AIFS architecture. These results demonstrate that the negative
 505 pre-activation field encodes valuable information regarding a proxy for saturation deficit. We
 506 acknowledge that these correlations are computed from a single 5-day forecast, which limits
 507 the temporal sampling. However, the analysis is performed on a Gaussian reduced N320 grid,

508 such that each 6-hourly forecast field contains more than 500,000 spatial evaluation points.
 509 Although based on one forecast initialization, the large number of grid-point samples per lead
 510 time provides a substantial statistical basis for examining the internal behaviour of the model.
 511

512 4.1.2 Optimization Geometry at the Zero-Precipitation Boundary

513 Having established that the negative pre-activation space encodes physically meaningful information,
 514 we now turn to understanding why constraining it during training improves forecast skill for
 515 light precipitation. The mechanism can be understood by examining how the Mean Squared
 516 Error (MSE) interacts with model outputs in the vicinity of the zero-precipitation boundary
 517 for a non-bounded model:

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

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

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

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

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

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

525 Because non-physical negative dry predictions (Scenario A) and genuine drizzle underpredictions
 526 (Scenario B) produce identical upward gradients, the optimizer receives an ambiguous training
 527 signal in the vicinity of zero. The loss provides no information about why the correction is
 528 required — whether it reflects a physical regime transition (dry \rightarrow drizzle) or merely a violation
 529 of the non-negativity constraint. One might expect the model to self-organize by learning to
 530 place dry predictions in a compact negative range — say, around -0.1 mm — thereby avoiding
 531 interference with the light-rain regime. However, this equilibrium is dynamically unstable
 532 under MSE. A dry prediction at -0.1 mm receives the same upward gradient as a genuine
 533 drizzle underprediction, so stochastic gradient updates continually push dry samples toward
 534 and across zero. As a result, no stable attractor can form in the negative space.

535 Importantly, the instability is locally asymmetric around $tp = 0$. For small $tp = \epsilon$ with
 536 $|\epsilon| \ll 1$,

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

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

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

539

Then

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

540

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

541

542

543

544

545

546

547

548

549

550

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

551

552

553

554

555

556

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

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

Figure 17 allows the gradient-ambiguity argument to be verified quantitatively. The three panels form a closed chain of evidence. The non-bounded model produces dry or negative outputs at only $\sim 10\%$ of grid points, compared to $\sim 30\%$ for the bounded model. The top-right panel shows that the non-bounded model’s light-precipitation frequency is inflated by almost exactly the same ~ 20 percentage points. The non-bounded model is not detecting more drizzle; it is misclassifying displaced dry events as light rain. The bottom panel reveals the mechanism predicted by the expected-gradient analysis. The non-bounded model produces a narrow spike of density straddling zero, within which the dry and drizzle regimes are superimposed and statistically indistinguishable. The distribution is tightly concentrated near zero but exhibits a slight positive skew, consistent with the theoretical result that the local stationary point of the expected MSE gradient lies at a strictly positive value. In other words, the model attempts to encode dry states in the neighbourhood of zero, yet the systematic upward drift induced by $\mathbb{E}[\partial \mathcal{L} / \partial tp] < 0$ for $tp < \mu$ prevents zero from acting as a stable attractor. The consequence is a persistent displacement of dry samples into weakly positive values, producing the observed excess of light precipitation.

572

573

574

575

576

577

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

578

579

580

581

582

583

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

584

585

The present analysis isolates the role of non-negativity enforcement and demonstrates that it addresses a fundamental gradient asymmetry at the zero-precipitation boundary. This

586 [mechanism operates at the level of local optimization dynamics and provides a distinct, mechanistically](#)
587 [interpretable pathway for drizzle reduction. While Sha et al. \(2025a\) demonstrate effectiveness](#)
588 [of combining non-negativity with global conservation constraints, our analysis suggests that](#)
589 [non-negativity merits investigation as an independent design element. The relative contributions](#)
590 [of boundary enforcement versus conservation-based regularization, and their potential architecture](#)
591 [dependence, remain important questions for future work.](#)

592 4.2 Case Studies

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

599 4.2.1 Storm Éowyn

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

605 Figure 18 shows the AIFS and IFS forecasts at decreasing lead times for the affected area
606 versus the corresponding IFS short-range forecast. The AIFS delivers an accurate forecast of
607 snowfall for this extremely rare event. This showcases the ability of the model to accurately
608 interpret meteorological patterns and forecast physically plausible events, even if they are far
609 from the training data. The AIFS predicted the event with a lead time of 10 days, earlier than
610 the IFS.

611 4.2.2 Tropical Low and extreme precipitation totals in Queensland Aus- 612 tralia

613 Starting in late January 2025, a slow-moving summer storm brought exceptional rainfall along
614 the northeastern coast of Queensland, Australia. Within a week, rainfall accumulation to-
615 talled more than 1000 millimetres in some areas, according to the Bureau of Meteorology as
616 reported in NASA Earth Observatory (2025). The city of Townsville saw the equivalent of
617 six months of rain in just three days and the largest weekly rainfall total was measured at
618 a gauge in the Cardwell Range, southwest of Tully, where nearly 1700mm fell (NASA Earth
619 Observatory (2025), Bureau of Meteorology measurements). Figure 19 compares forecasts
620 from AIFS and IFS against the IMERG Huffman et al. (2023) final product for the period
621 01/02/2025–03/02/2025. Both model forecasts were initialized on 30/01/2025, two days prior
622 to the event. The Cardwell Range is indicated by a black star, and the city of Townsville by a
623 cyan star. Both IFS and AIFS successfully captured the event, with 24-hour rainfall accumu-
624 lations exceeding 300 mm in some regions. However, the AIFS forecast exhibits a somewhat
625 persistent signal in the 5-day lead time, predicting very high rainfall totals near the Cardwell
626 Range. This highlights that, despite AIFS’s tendency toward excessive spatial smoothing, it
627 remains capable of accurately forecasting extreme events at medium range.

5 Discussion and conclusion

The revised AIFS version (1.1.0) presented here improves upon the pre-operational release through a revised training regime with more data, new forecast variables, improved stratospheric loss weights, and a bounding strategy that enforces physical constraints on the output variables. Overall, this leads to improvements of around 4–6 % across all variables, lead times, and pressure levels. The largest improvements, up to 12% gains in normalized difference in the short range, are observed in total precipitation forecasting, which benefits from the newly introduced bounding. We showed that this has a significant impact on the prediction of no rain and light precipitation. The model displays good forecast performance for out-of-training-sample case studies, accurately capturing extreme precipitation and snowfall events.

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

~~The bounding strategy implemented also plays a crucial role, especially for precipitation forecasting. Hard-constraining model outputs increases the physical realism of forecast fields and the light precipitation forecast skill. This improvement is attributed to a shift in the forecast space, where bounding the output facilitates the prediction of~~ “Recent global AI forecasting systems, including GraphCast, Pangu-Weather, FuXi, and CREDIT, have reported persistent challenges in representing light precipitation. Positive frequency bias in the drizzle regime appears to be a recurring feature across models trained with symmetric regression losses on strictly non-negative, intermittent variables. Although these systems differ substantially in backbone architecture, from graph neural networks to transformer-based designs and modular physically constrained frameworks, the drizzle problem appears largely independent of architecture. Instead, it is closely tied to how precipitation is parameterized and constrained during training. Physical constraint methodologies offer multiple pathways for mitigating precipitation biases. Global conservation schemes may reduce drizzle indirectly by regulating total moisture budgets. The present analysis suggests that non-negativity enforcement addresses a more fundamental issue: the local gradient asymmetry at the zero boundary and the superposition of dry and wet states around zero. By introducing a hard geometric boundary at zero, the optimization landscape is reshaped such that dry and wet regimes become separable. This mechanism operates independently of large-scale conservation principles and may therefore represent a structural requirement for stable training of intermittent variables under MSE. Alternative activation functions such as LeakyReLU, which scale negative inputs by a small factor α (typically 0.01), would permit gradient flow in the negative space while heavily attenuating the loss contribution from dry predictions (by a factor of α^2). We expect that similar regime separation would still emerge, since the cost of placing dry states deep in the negative space becomes negligible. The main practical difference is that LeakyReLU produces non-physical slightly negative output values at inference, requiring post-processing clipping. More broadly, alternative formulations that explicitly decouple the dry and wet regimes during training, such as asymmetric loss functions or dedicated classification heads for the no-rain”. We hypothesise that the constraint enables the model to treat the negative space as likelihood of “no-rain” conditions, thereby facilitating the learning of the zero tp forecast. Other bounding functions may be used, and we plan to explore these in future work. In particular, we aim to investigate LeakyReLU-based approaches, which allow for weight updates with changes in the negative space, something that standard ReLU functions do not permitstate, represent

677 promising directions for future work.

678 The bounding strategy presented here enforces physical realizability, non-negativity, boundedness,
679 and inter-variable consistency, but does not impose global conservation of mass or energy. For
680 the medium-range timescales considered in this work (up to 10 days), we expect conservation
681 violations to remain small relative to forecast errors dominated by chaotic error growth, though
682 a systematic quantification of mass and energy drift over extended AIFS integrations remains
683 to be carried out.

684 ~~Rollout fine-tuning also emerges as a key factor shaping the forecasting skill of the model,~~
685 ~~particularly through its influence on the smoothing of outputs. While smoothing is already~~
686 ~~present in the pre-trained model, rollout fine-tuning enhances this behaviour. This reflects~~
687 ~~the model's adaptation to the inherent forecast uncertainty for longer lead times. emerges~~
688 ~~as an important factor shaping forecast behaviour, including the degree of spatial smoothing~~
689 ~~in the outputs. As the model is exposed to lead times up to 72h, the minimization of the~~
690 ~~trained on extended lead times and optimised using a mean squared error inevitably results in~~
691 ~~an enhanced blurring of the fields. The impact of training hyperparameters on the smoothing~~
692 ~~characteristics of the output remains to be further explored. Limited testing has shown that~~
693 ~~factors objective, some degree of smoothing is expected. Training hyperparameters such as~~
694 ~~learning rate scheduling, number of steps and the rollout strategy all have an influence on~~
695 ~~the intensity of blurring in the fields. These findings highlight an important design trade-off~~
696 ~~in training deterministic AI forecasting models: between forecast realism and optimization~~
697 ~~of MSE-based verification scores. Forecast realism refers to how physically plausible and~~
698 ~~meteorologically coherent AI-generated forecasts are, here specifically in terms of their spectral~~
699 ~~characteristics (e.g. power spectra across spatial scales). One way to evaluate the physical~~
700 ~~realism of the resulting forecasts is to assess whether their spectral signature resembles that~~
701 ~~observed in the analysis. While aggressive optimisation steps, and rollout fine-tuning strategies~~
702 ~~(such as the one used in Bodnar et al. (2025)) can significantly boost headline scores, here we~~
703 ~~have chosen an approach that maintains a subjective compromise between forecast realism and~~
704 ~~forecast skill measured by RMSE configuration can influence this behaviour and warrant further~~
705 ~~systematic investigation. In the present study, the training configuration, including a maximum~~
706 ~~rollout length of 12, was retained from the previous AIFS version to ensure consistency. The~~
707 ~~resulting Z500 power spectra (Figure 6) are broadly comparable to those of the previous model~~
708 ~~across scales, including the 500 km range, with slightly improved agreement with the analysis at~~
709 ~~longer lead times. Importantly, these comparable spectral characteristics are achieved alongside~~
710 ~~overall improvements in RMSE-based skill (Figure 7). This indicates that the skill gains are~~
711 ~~not obtained at the expense of degraded spatial variability. While more aggressive rollout~~
712 ~~strategies may further optimise headline verification scores, understanding their impact on~~
713 ~~spatial characteristics remains an important area for future work.~~

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

727 (ORAS6) (Zuo et al., 2024).

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

730 The AIFS development has now transitioned to the new Anemoi framework (Lang et al.,
731 2024a; Nipen et al., 2024; Wijnands et al., 2025). Anemoi provides tools for the whole data-
732 driven modelling workflow, from the generation of training datasets, to scalable probabilistic
733 training (Lang et al., 2024b) and running real-time inference with such models. Anemoi also
734 allows for the cataloguing and archiving of model and data checkpoints to ensure reproducibility
735 and traceability of training and inference runs and ensure that any models developed within
736 this framework have a clear lineage. The Anemoi framework is now being used by an increasing
737 number of Member States of ECMWF and collaborating organisations supported by ECMWF.

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

748 6 Code and Model Availability

749 AIFS version 1.1.0 was fully trained using the Anemoi framework [https://github.com/ecmwf/](https://github.com/ecmwf/anemoi)
750 [anemoi](https://github.com/ecmwf/anemoi). The frozen versions of the Anemoi modules used for training, together with the config-
751 uration files and the trained model checkpoint, are available in the permanent archive European
752 Centre for Medium-Range Weather Forecasts (2025) under DOI [https://doi.org/10.5281/](https://doi.org/10.5281/zenodo.17349820)
753 [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
754 Hugging Face <https://huggingface.co/ecmwf/aifs-single-1.1> under a Creative Commons
755 Attribution 4.0 International (CC BY 4.0) licence and DOI [https://doi.org/10.57967/hf/](https://doi.org/10.57967/hf/6415)
756 [6415](https://doi.org/10.57967/hf/6415) (ECMWF, 2025a).

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

763 Author Contributions

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

769 *Equal Contribution

7 Acknowledgements

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

References

- Blanka Balogh, David Saint-Martin, and Olivier Geoffroy. Online test of a neural network deep convection parameterization in arp-gem1, 2024. URL <https://arxiv.org/abs/2410.21920>.
- Zied Ben Bouallègue, Mariana C A Clare, Linus Magnusson, Estibaliz Gascón, Michael Maier-Gerber, Martin Janoušek, Mark Rodwell, Florian Pinault, Jesper S Dramsch, Simon T K Lang, Baudouin Raoult, Florence Rabier, Matthieu Chevallier, Irina Sandu, Peter Dueben, Matthew Chantry, and Florian Pappenberger. The rise of data-driven weather forecasting: A first statistical assessment of machine learning-based weather forecasts in an operational-like context. *Bulletin of the American Meteorological Society*, 2024. ISSN 1520-0477. doi: [doi: doi.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>.
- K. Bi, L. Xie, H. Zhang, et al. Accurate medium-range global weather forecasting with 3D neural networks. *Nature*, 619:533–538, 2023. doi: [10.1038/s41586-023-06185-3](https://doi.org/10.1038/s41586-023-06185-3).
- Cristian Bodnar, Wessel P Bruinsma, Ana Lucic, Megan Stanley, Anna Allen, Johannes Brandstetter, Patrick Garvan, Maik Riechert, Jonathan A Weyn, Haiyu Dong, et al. A foundation model for the earth system. *Nature*, 641:1180–1187, 2025. doi: [10.1038/s41586-025-09005-y](https://doi.org/10.1038/s41586-025-09005-y).
- Massimo Bonavita. On some limitations of current machine learning weather prediction models. *Geophysical Research Letters*, 51(12):e2023GL107377, 2024. doi: <https://doi.org/10.1029/2023GL107377>.
- Boris Bonev, Thorsten Kurth, Ankur Mahesh, Mauro Bisson, Jean Kossaifi, Karthik Kashinath, Anima Anandkumar, William D. Collins, Michael S. Pritchard, and Alexander Keller. Four-castnet 3: A geometric approach to probabilistic machine-learning weather forecasting at scale, 2025. URL <https://arxiv.org/abs/2507.12144>.
- Noah D. Brenowitz, Yair Cohen, Jaideep Pathak, Ankur Mahesh, Boris Bonev, Thorsten Kurth, Dale R. Durran, Peter Harrington, and Michael S. Pritchard. A practical probabilistic benchmark for ai weather models. *Geophysical Research Letters*, 52(7), April 2025. ISSN 1944-8007. doi: [10.1029/2024gl113656](https://doi.org/10.1029/2024gl113656). URL <http://dx.doi.org/10.1029/2024GL113656>.
- Lei Chen, Xiaohui Zhong, Feng Zhang, Yuan Cheng, Yinghui Xu, Yuan Qi, and Hao Li. FuXi: a cascade machine learning forecasting system for 15-day global weather forecast. *npj Climate and Atmospheric Science*, 6(1), November 2023. ISSN 2397-3722. doi: [10.1038/s41612-023-00512-1](https://doi.org/10.1038/s41612-023-00512-1). URL <http://dx.doi.org/10.1038/s41612-023-00512-1>.
- ECMWF. aifs-single-1.1 (revision 7976552), 2025a. URL <https://huggingface.co/ecmwf/aifs-single-1.1>.

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

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

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

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

822 H. Hersbach, B. Bell, P. Berrisford, et al. The ERA5 global reanalysis. *QJ R Meteorol Soc*,
823 146:1999–2049, 2020. doi: 10.1002/qj.3803.

824 George J. Huffman, Erich F. Stocker, David T. Bolvin, Eric J. Nelkin, and Jackson Tan. GPM
825 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.
826

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

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

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

834 Remi Lam, Alvaro Sanchez-Gonzalez, Matthew Willson, Peter Wirsberger, Meire Fortu-
835 nato, Ferran Alet, Suman Ravuri, Timo Ewalds, Zach Eaton-Rosen, Weihua Hu, Alexan-
836 der Merose, Stephan Hoyer, George Holland, Oriol Vinyals, Jacklynn Stott, Alexander
837 Pritzel, Shakir Mohamed, and Peter Battaglia. Learning skillful medium-range global
838 weather forecasting. *Science*, 382(6677):1416–1421, December 2023. ISSN 1095-9203. doi:
839 10.1126/science.adi2336. URL <http://dx.doi.org/10.1126/science.adi2336>.

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

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

851 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International*
852 *Conference on Learning Representations*, 2019. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=Bkg6RiCqY7)
853 [Bkg6RiCqY7](https://openreview.net/forum?id=Bkg6RiCqY7).

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

857 J. Muñoz Sabater, E. Dutra, A. Agustí-Panareda, C. Albergel, G. Arduini, G. Balsamo,
858 S. Bousssetta, M. Choulga, S. Harrigan, H. Hersbach, B. Martens, D. G. Miralles, M. Piles,
859 N. J. Rodríguez-Fernández, E. Zsoter, C. Buontempo, and J.-N. Thépaut. Era5-land: a
860 state-of-the-art global reanalysis dataset for land applications. *Earth System Science Data*,
861 13(9):4349–4383, 2021. doi: 10.5194/essd-13-4349-2021. URL <https://essd.copernicus.org/articles/13/4349/2021/>.

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

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

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

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

879 Stephan Rasp, Stephan Hoyer, Alexander Merose, Ian Langmore, Peter Battaglia, Tyler
880 Russell, Alvaro Sanchez-Gonzalez, Vivian Yang, Rob Carver, Shreya Agrawal, Matthew
881 Chantry, Zied Ben Bouallegue, Peter Dueben, Carla Bromberg, Jared Sisk, Luke Barrington,
882 Aaron Bell, and Fei Sha. Weatherbench 2: A benchmark for the next generation
883 of data-driven global weather models. *Journal of Advances in Modeling Earth Systems*,
884 16(6):e2023MS004019, 2024. doi: <https://doi.org/10.1029/2023MS004019>. URL <https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2023MS004019>. e2023MS004019
885 2023MS004019.

887 Mark J. Rodwell, David S. Richardson, Tim D. Hewson, and Thomas Haiden. A new equitable
888 score suitable for verifying precipitation in numerical weather prediction. *Quarterly Journal*
889 *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>.

891 John S. Schreck, Yingkai Sha, William Chapman, Dhamma Kimpara, Judith Berner, Seth
892 McGinnis, Arnold Kazadi, Negin Sobhani, Ben Kirk, Charlie Becker, Gabrielle Gantos,
893 and David John Gagne II. Community research earth digital intelligence twin: a scalable
894 framework for ai-driven earth system modeling. *npj Climate and Atmospheric Science*, 8(1),
895 June 2025. ISSN 2397-3722. doi: 10.1038/s41612-025-01125-6. URL <http://dx.doi.org/10.1038/s41612-025-01125-6>.

897 Yingkai Sha, John S. Schreck, William Chapman, and David John Gagne II. Investi-
898 gating the use of terrain-following coordinates in ai-driven precipitation forecasts. *Geo-*
899 *physical Research Letters*, 52(20):e2025GL118478, 2025a. doi: [https://doi.org/10.1029/](https://doi.org/10.1029/2025GL118478)
900 [2025GL118478](https://doi.org/10.1029/2025GL118478). URL [https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/](https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2025GL118478)
901 [2025GL118478](https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2025GL118478). e2025GL118478 2025GL118478.

902 Yingkai Sha, John S. Schreck, William Chapman, and David John Gagne II. Improving ai
903 weather prediction models using global mass and energy conservation schemes. *Journal of*
904 *Advances in Modeling Earth Systems*, 17(11):e2025MS005138, 2025b. doi: [https://doi.org/](https://doi.org/10.1029/2025MS005138)
905 [10.1029/2025MS005138](https://doi.org/10.1029/2025MS005138). URL [https://agupubs.onlinelibrary.wiley.com/doi/abs/10.](https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2025MS005138)
906 [1029/2025MS005138](https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2025MS005138). e2025MS005138 2025MS005138.

907 Akshay Subramaniam, Dale Durran, David Pruitt, Nathaniel Cresswell-Clay, and William Yik.
908 Imposing the fundamental dynamical constraint of hydrostatic balance to improve global ml
909 weather prediction, 2025. URL <https://arxiv.org/abs/2506.08285>.

910 Haley Thiem and Nicole Collins. Historic January 2025 snowstorm in the south-
911 ern US, 2025. URL [https://www.climate.gov/news-features/event-tracker/](https://www.climate.gov/news-features/event-tracker/historic-january-2025-snowstorm-southern-us)
912 [historic-january-2025-snowstorm-southern-us](https://www.climate.gov/news-features/event-tracker/historic-january-2025-snowstorm-southern-us).

913 N. P. Wedi. Increasing the horizontal resolution in numerical weather prediction and climate
914 simulations: illusion or panacea? *Philosophical Transactions of the Royal Society A*, 372,
915 2014. doi: 10.1098/rsta.2013.0289.

916 Jasper S. Wijnands, Michiel Van Ginderachter, Bastien François, Sophie Buurman, Piet Term-
917 nia, and Dieter Van den Bleeken. A comparison of stretched-grid and limited-area modelling
918 for data-driven regional weather forecasting, 2025. URL [https://arxiv.org/abs/2507.](https://arxiv.org/abs/2507.18378)
919 [18378](https://arxiv.org/abs/2507.18378).

920 Daniel S Wilks. *Statistical methods in the atmospheric sciences*. Elsevier Science Publishing,
921 Philadelphia, PA, 4 edition, June 2019.

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

(a) Convective precipitation (b) Total precipitation (c) Difference (cp-tp) 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.

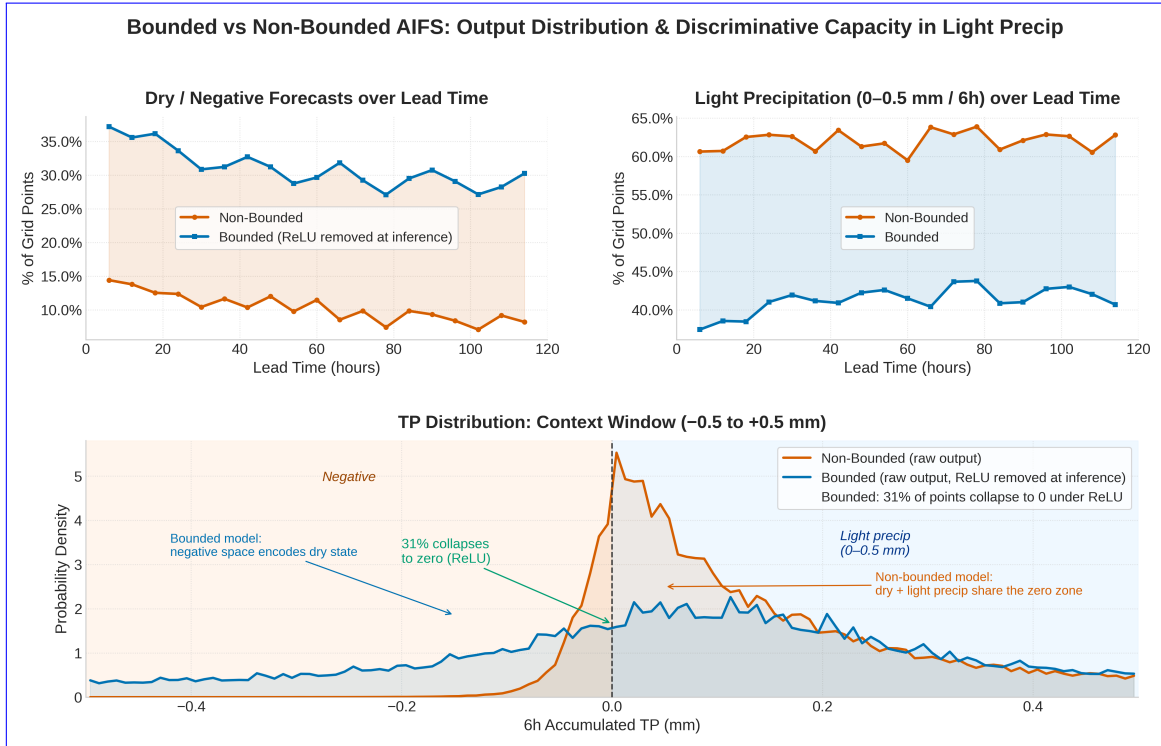


Figure 17: Output distribution and discriminative capacity in the light-precipitation regime for bounded and non-bounded AIFS. The bounded model's ReLU is removed at inference to expose raw pre-activations. (Top left) The non-bounded model produces dry or negative outputs at only ~10% of grid points versus ~30% for the bounded model. A persistent 20-percentage-point gap across all lead times. (Top right) The non-bounded model assigns ~60% of grid points to the light-precipitation bin (0-0.5 mm / 6h) versus ~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.

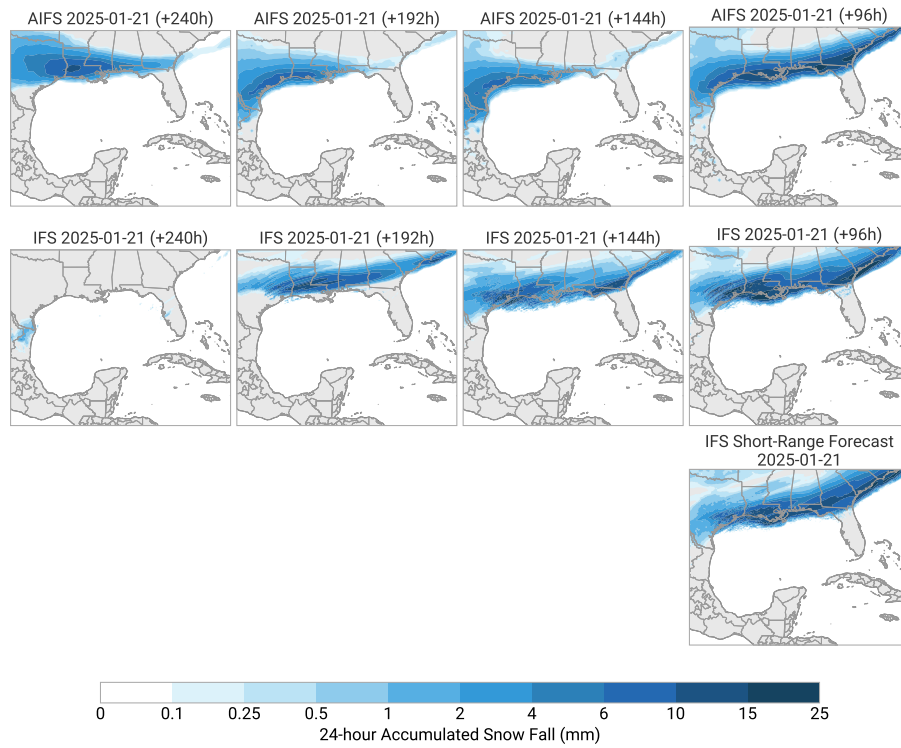


Figure 18: Snowfall forecasts for AIFS (top row) and IFS (middle row) over the Gulf Coast of America at 10, 8, 6 and 4 day lead times from left to right respectively, against IFS short-range forecasts for the snowfall event (bottom row). The figure shows how the snowfall event was forecast accurately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

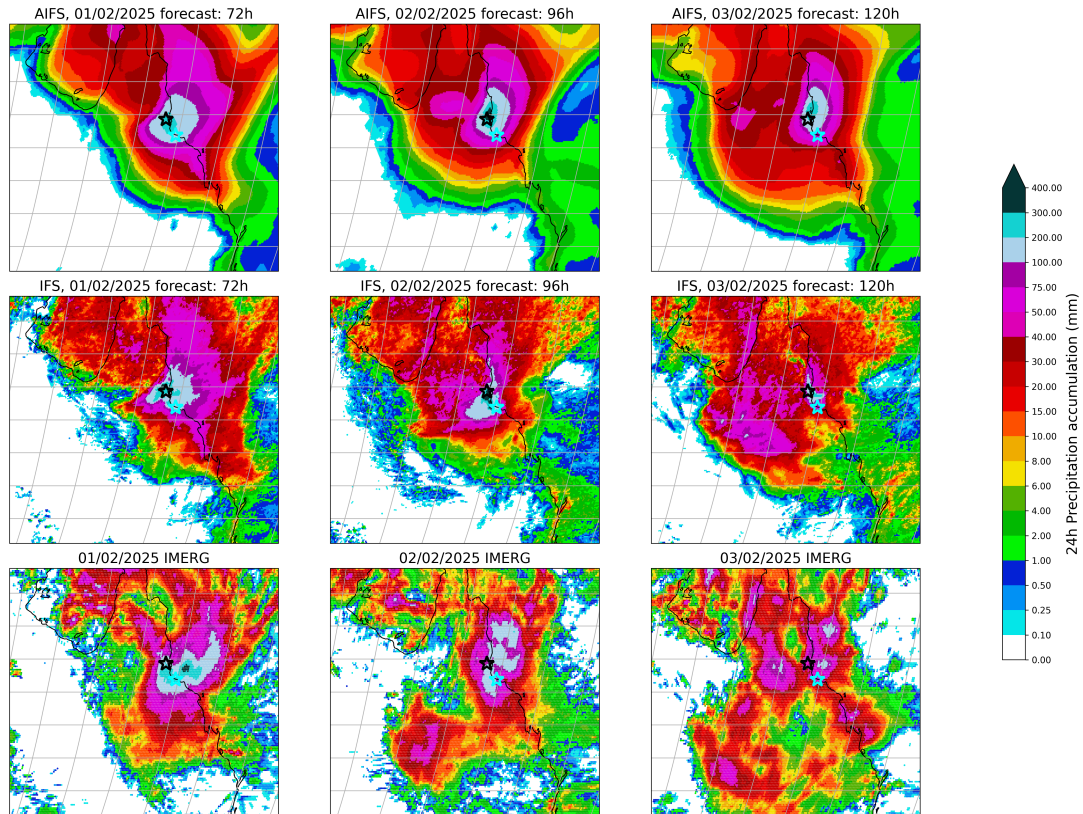


Figure 19: 24-hour accumulated precipitation forecasts from the AIFS (top row) and IFS (middle row) models, compared with IMERG observational data (bottom row) over northeastern Queensland for 01/02/2025 to 03/02/2025. Forecasts are initialised on 30/01/2025. The black star marks the Cardwell Range, 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.