

1 **Response to Editor**

2 **Dear authors,**

3 **Thank you for submitting the revised version of your paper to GMD.**

4 **As you can see, one of the reviewers has provided a thorough and constructive review of your**
5 **revised manuscript. I agree with the reviewer that it is necessary to evaluate the short-term skill**
6 **of the model, rather than relying on climatological precipitation, and also use a statistically**
7 **independent dataset. Both issues need to be addressed.**

8 **Regarding code and data availability, I also agree that although GPM and ERA5 can largely**
9 **be considered standard datasets in atmospheric sciences, it would be appropriate to make the**
10 **specific dataset you used available via a Zenodo archive in order to make the results of the paper**
11 **fully reproducible.**

12 **Kind regards,**

13 **Emmanouil Flaounas**

14

15 **Dear Editor and Reviewers,**

16 **Thank you for your time, review, and consideration of our manuscript. We have revised the paper**
17 **in response to your comments. In particular, we re-evaluated the offline training performance using a**
18 **fully independent GSRM dataset that is outside the 80-day simulation samples. We also added an**
19 **assessment of the short-term forecast skill of GRIST-MPS versus GRIST-CPS, which is provided in**
20 **the Supplementary Materials. The specific dataset we used has already been stored at Zenodo, and is**
21 **now explicitly stated in the code and data availability.**

22 **All other comments have been addressed point by point in the response.**

23

24 **All the best**

25 **Yi Zhang (on behalf of the authors)**

26

27 Response to Reviewers' Comments

28 Response to Reviewer #1:

29 We thank this reviewer for the insightful comments and detailed suggestions on how to improve
30 the manuscript. The manuscript has been further improved based on your comments. Below is the
31 item-by-item reply to your questions and suggestions, the texts with normal font are your original
32 comments, the texts with blue font are our responses and the texts with italics are the revised content
33 of the manuscript.

34 Line comments:

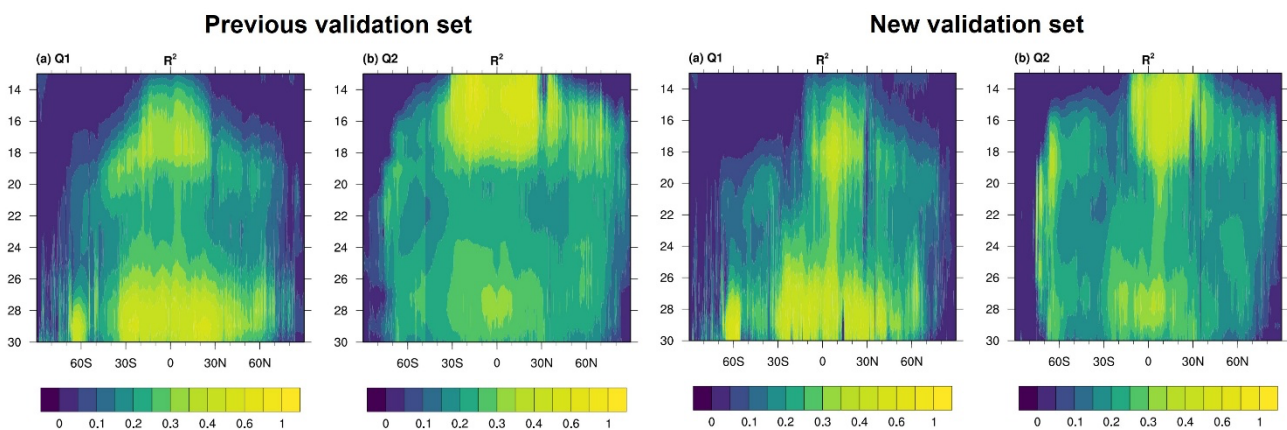
35 **L118: This split does not yield an independent validation dataset. Neighboring time points in**
36 **your dataset are incredibly correlated, especially because they are derived using linear**
37 **interpolation from a coarser time grid. The linear interpolation aspect means you might as well**
38 **have used the training dataset itself as the validation dataset. This is not a reliable assessment of**
39 **out-of-sample performance. Assuming Figure 2 is based on validation data, it must be**
40 **reevaluated using a statistically independent validation dataset, as this is your main measure of**
41 **process skill.**

42 Reply: We appreciate the reviewer's suggestion. Accordingly, we introduced a completely independent
43 10-day GSRM dataset as a new validation set to re-evaluate the model's out-of-sample performance.
44 We have updated Figure 2 in the main text accordingly.

45 As also shown in Figure A1, the overall patterns remain consistent between the two validation sets.
46 While the R^2 values from the new validation set (right panel) are slightly lower in certain regions
47 compared with those from the previous validation set (left panel). This indicates that the model retains
48 a reasonable degree of reliability on out-of-sample data.

49 Corresponding revisions have also been incorporated into the manuscript. (L119-122).

50 *“Moreover, we have further conducted a fully independent 10-day GSRM experiment (date: 14–*
51 *23th, July 2008) for final validation to assess the model's out-of-sample performance (see Fig. 2 for*
52 *details).”*



53
54 Figure A1. Offline skill of the coefficient of determination (R^2) for Q_1 and Q_2 , as functions of
55 latitude and model level. Left panel: the previous Figure 2, Right panel: as in Figure 2 but for new
56 independent 10-day GSRM validation set.

57 **L166: Remove “per epoch” from “weight decay of 10^{-6} per epoch” . Adam weight decay**
58 **is applied on each gradient update, not just once per epoch.**

59 Reply: Thank you for your comments. Done.

60 **L167: Please define “superior performance”. What metric on what data?**

61 Reply: We have now added the performance and data descriptions to the revised manuscript (L167-
62 168).

63 *“The mean absolute error (MAE) loss was selected over the mean squared error (MSE) loss as the*
64 *loss function, as it demonstrated superior performance with lower biases on validation sets during*
65 *initial training phases.”*

66 **L184: What “offline performance benchmarks” - what metric(s) on what data?**

67 Reply: Offline performance benchmarks refer to the metrics of $MSE < 1 \times 10^{-4}$ and $R^2 > 0.3$ achieved
68 on the validation sets as described above. This content has also been incorporated into the
69 corresponding paragraph. (L184-186)

70 *“The final selection of our optimal MPS is based on a dual evaluation: satisfying offline*
71 *performance benchmarks ($MSE < 1 \times 10^{-4}$ and $R^2 > 0.3$ on validation sets) and demonstrating stability*
72 *in online integration, which must maintain stable online integration for more than 3 months.”*

73 **L186: “can integrate stable” -> “can integrate stably”**

74 Reply: We have revised it to “demonstrated stable integration”.

75 **L186: Please define “better performance”, what metric on what data? Is this on validation data,**
76 **or is this a metric applied to the online simulation? To 6 years of simulation or a longer/shorter**
77 **period?**

78 Reply: "Better performance" refers to the ability to perform stable integration for more than 6 years,
79 with better RMSE values for online variables such as precipitation. Relevant details have been added
80 in the manuscript (L187-189).

81 *“Among the eight NNs, two demonstrated stable integration capabilities for over six years. We*
82 *selected the better-performing model between them as the optimal MPS, based on its improved RMSE*
83 *values for online variables such as precipitation (see Table 1).”*

84 **L237-244: Please comment in the manuscript on whether the model is unstable or simply**
85 **performs poorly when the near-surface or top-of-atmosphere levels use MPS instead of CPS.**

86 Reply: The model is unstable when the near-surface or top-of-atmosphere levels use MPS. The relevant
87 descriptions have been added in the manuscript. (L241-242)

88 *“This selective preservation, validated through sensitivity experiments, serves as a stability-*
89 *enhancing mechanism, the present MPS without this configuration will crash after running for a few*
90 *days.”*

91 **L267-273, Figure 4: In several areas away from the pacific ITCZ region, GRIST-MPS results**
92 **look clearly worse than the GRIST-CPS values. This is borne out in the pattern correlation**
93 **coefficients, which indicate 26% and 12% unexplained variance of the pattern respectively -**
94 **more than double in the MPS-based model. Specifically, there are large dry biases in the**
95 **northern and southern extratropical ocean winter regions, and a large moist bias over Africa in**
96 **both seasons. Provided this is an accurate assessment, please ensure these are remarked on, or**
97 **add to the paragraph on lines L288-290.**

98 Reply: Thank you for this remark. The content regarding the limitations of GRIST-MPS has now been
99 incorporated into the manuscript. (L287-289).

100 *“Meanwhile, some remaining biases persist in GRIST-MPS: compared to GRIST-CPS, its global*
101 *spatial pattern correlation coefficient is marginally lower (PCC: MPS: 0.86 vs.CPS: 0.94) on*
102 *Summer. For it exhibits substantial dry biases over the northern and southern subtropical oceans in*
103 *winter, alongside a consistent moist bias over Africa in both seasons.”*

104 **You may want to separately remark on the pattern correlation coefficient confined to the pacific**
105 **ITCZ region, which I would expect to be better, reflective of the improved ITCZ strength and**
106 **width.**

107 Reply: We calculated the pattern correlation coefficients for the Pacific ITCZ region, the results show
108 that CPS yields better performance in summer (MPS:0.89 vs. CPS: 0.93) , whereas MPS exhibits a
109 more robust simulation capacity in winter (MPS:0.95, CPS:0.94). We attribute this to the fact that
110 although both the strength and width have been improved in both seasons, there remain certain biases
111 in the position of the ITCZ compared with GPM during summer. Such biases give rise to a double
112 penalty issue, which ultimately results in its correlation coefficients being lower than those of CPS.
113 For this reason, the correlation coefficients for the Pacific ITCZ region are not presented in this section;
114 instead, the ITCZ and spatial correlation are introduced separately.

115 **L268-269: Wodzicki and Rapp 2016 do not define ITCZ strength in this way. In section 2.3 they**
116 **specify that precipitation intensity should be used to measure ITCZ strength, but they take this**
117 **as a monthly mean value across grid cells in a specific ITCZ region. Please use established**
118 **measures for ITCZ strength and width.**

119 Reply: The strength and width of the ITCZ have been recalculated using well-defined methods, with
120 appropriate citations provided for each (L271-276).

121 *“Following established conventions, we quantify the width of the ITCZ as the latitudinal distance*
122 *between its northern and southern boundaries, applying a 5-point smoother to the data prior to*
123 *calculation (Wodzicki and Rapp, 2016). And the northern (southern) boundary is identified by moving*
124 *equatorward from higher latitudes and locating the first grid cell where the precipitation in the*
125 *adjacent grid cell to the north (south) falls below 2.5 mm/day. The strength of the ITCZ is defined as*
126 *the area-weighted mean precipitation within the region bounded by these northern and southern*
127 *boundaries (Wang et al., 2023).”*

128 **L270: The cited Byrne et al. 2018 paper is a review paper which does not define ITCZ width.**
129 **They reference one paper which defines ITCZ extent using a precipitation minimum - Wodzicki**
130 **and Rapp 2016 (<https://doi.org/10.1002/2015JD024458>), which you cited for precipitation**
131 **strength. This paper reviewed a range of rainfall thresholds (1, 2.5, and 5.0 mm/day) and found**
132 **2.5 mm/day to be most consistent across datasets, and used a 5-point smoother prior to taking**
133 **the threshold. Please use this methodology for measures of ITCZ width and cite this paper, or**
134 **cite another paper supporting a different methodology.**

135 Reply: We appreciate your suggestion. We have now applied the method you suggested, and provided
136 a clear description of it in the paper, as detailed in the previous response.

137 **Summary and outlook:**

138 **It is hard for me to agree with publishing the characterization of model skill given in this**
139 **summary. While it is true that the model more accurately represents the Pacific ITCZ rainband**
140 **(thank you for improvements to the manuscript which quantify this), significant precipitation**
141 **biases are produced in most other regions with strong precipitation features, worse than the CPS-**
142 **based model. This is borne out in the pattern correlation coefficients. These should be given**
143 **proportionate weight in the assessment. The relevance of this work for future ML schemes could**
144 **be discussed in more detail.**

145 Reply: Thank you for your comments. The limitations of GRIST-MPS have been clearly indicated in
146 the manuscript(L287-289), and the discussion on the relevance of this study for future ML schemes
147 has also been reorganized accordingly(L392-410).

148 “Meanwhile, some remaining biases persist in GRIST-MPS: compared to GRIST-CPS, its global
149 spatial pattern correlation coefficient is marginally lower (PCC: MPS: 0.86 vs. CPS: 0.94). For it
150 exhibits substantial dry biases over the northern and southern subtropical oceans in winter, alongside
151 a consistent moist bias over Africa in both seasons.”

152 “**Prospective Relevance.** As a proof-of-concept study, this work provides useful reference for
153 future ML-based efforts by demonstrating that a column-based ML physics module, trained on GSRM
154 multiscale modeling data over a limited time window, can produce realistic online free-running climate
155 simulations and generalize to unseen periods. Beyond the choice of ML architecture, several strategies
156 proved essential. First, spatial random sampling mitigated spatial overfitting arising from an
157 imbalanced data distribution, improving both accuracy and stability. Second, interpolating the data
158 from 1-hour to 20-minute resolution tripled the effective sample size, imparted a regularization effect,
159 and aligned the training data with the model time step, further improving online performance. Notably,
160 the model still requires a small number of near-surface CPS tendencies (primarily for the boundary
161 layer) to maintain online stability. This likely reflects a limitation of the coarse-grained dataset, which
162 supports only a coarse vertical resolution in the lower troposphere (spacing >200 m below 850 hPa).
163 In addition, vertical eddy transport associated with boundary-layer turbulence may require higher-
164 resolution training data to be represented effectively.”

165 “An important implication is that representing physics tendencies as residual terms of grid-scale
166 variables is a promising route for diagnosing training targets for ML-based physics. This approach
167 offers several potential advantages that merit further exploration. For example, a common high-
168 resolution dataset can be incrementally coarse-grained to multiple coarser resolutions, enabling
169 datasets at different resolutions to share consistent large-scale information. Moreover, the large-scale
170 states can be further constrained (e.g., via short-period simulations and/or nudging), which may
171 improve the quality of the diagnosed tendencies, and/or generate additional constraint-related
172 tendencies that can be learned separately. However, unlike physics tendencies directly extracted from
173 a host model, this reconstructed-tendency approach generally requires dedicated training strategies
174 and careful coupling procedures to ensure stable and accurate online integration.”

175 **Code and data availability: The editor should comment on this, but is it sufficiently reproducible**
176 **to provide a simple web link to the top level page for GPM and ERA5 data? Websites are subject**
177 **to change, the process for downloading this data is quite involved, and at least in the case of**
178 **ERA5 there are several data products to choose from.**

179 Reply: The specific ERA5 and GPM datasets used in this study are also available in the linked code
180 repository, within the compressed archive [input_plot.tar.gz](#). This has been now explicitly mentioned
181 in the data availability.

182 “The specific GPM and ERA5 datasets used in this work were archived in the file [input_plot.tar.gz](#)
183 on the Zenodo repository provided.”

184 **Overall review:**

185 **It looks like a wash whether this model outperforms the CPS-based model on climate skill -**
186 **arguably the precipitation climatology is worse. I do not think this on its own is a reason to reject**
187 **this manuscript. However, it is important that the manuscript does not give a false sense of the**
188 **climate performance.**

189 Reply: Thank you for your valuable comments. In the original manuscript, we presented results
190 primarily emphasizing where the GRIST-MPS model outperformed the GRIST-CPS model. As you
191 noted, this argument is not fully conclusive based on the present simulation, but it does indicate a

192 positive signal of learning physics from GSRMs.

193 In the revised version, we have accordingly strengthened the analysis of precipitation climatology
194 and, more importantly, elaborated on the limitations of the GRIST-MPS model to provide a more
195 balanced and objective evaluation.

196 *“Meanwhile, some remaining biases persist in GRIST-MPS: Compared to GRIST-CPS, its global
197 spatial pattern correlation coefficient is marginally lower (PCC: MPS: 0.86 vs. CPS: 0.94). For it
198 exhibits substantial dry biases over the northern and southern subtropical oceans in winter, alongside
199 a consistent moist bias over Africa in both seasons.”*

200 **I do think it’s an achievement to produce ML-based column physics trained on only 80 days of
201 data which does anything remotely reasonable when run online, and has any accuracy on out-
202 of-sample process representations. It is likely the use of CPS near the surface and top of
203 atmosphere is crucial for this, as well as the use of interpolation to augment the training dataset
204 or regularize the model. If this is the main achievement of the paper, it is important to
205 demonstrate that these are key components by commenting on the performance or stability when
206 these features are omitted.**

207 Reply: We appreciate your recognition. The experiments you mentioned have been conducted and are
208 summarized in Table 3, with the exception of the test using the near-surface MPS, which led to rapid
209 model instability.

210 As shown in Table 3, both random sampling and temporal interpolation contribute significantly to
211 model performance: random sampling effectively improves stability, while temporal interpolation
212 primarily enhances accuracy, and also suggests better generalization ability in models trained with this
213 approach. Corresponding discussions of these findings are included in the concluding section of the
214 paper.

215 *“**Prospective Relevance.** As a proof-of-concept study, this work provides useful reference for future
216 ML-based efforts by demonstrating that a column-based ML physics module, trained on GSRM
217 multiscale modeling data over a limited time window, can produce realistic online free-running climate
218 simulations and generalize to unseen periods. Beyond the choice of ML architecture, several strategies
219 proved essential. First, spatial random sampling mitigated spatial overfitting arising from an
220 imbalanced data distribution, improving both accuracy and stability. Second, interpolating the data
221 from 1-hour to 20-minute resolution tripled the effective sample size, imparted a regularization effect,
222 and aligned the training data with the model time step, further improving online performance. Notably,
223 the model still requires a small number of near-surface CPS tendencies (primarily for the boundary
224 layer) to maintain online stability. This likely reflects a limitation of the coarse-grained dataset, which
225 supports only a coarse vertical resolution in the lower troposphere (spacing >200 m below 850 hPa).
226 In addition, vertical eddy transport associated with boundary-layer turbulence may require higher-
227 resolution training data to be represented effectively.”*

228 **The use of highly correlated validation data is still a major issue in the manuscript which must
229 be addressed. Figure 2 must be produced using samples which are not correlated to the training
230 dataset. If this figure does not show skill when evaluated on independent examples (which is a
231 possible result), it would not be reasonable to publish this manuscript. Ideally uncorrelated out-
232 of-sample validation would be completed throughout training process, to assess whether the
233 model is overfitting and at what point during training. This is not being done with the described
234 train validation split.**

235 Reply: As replied in our earlier response, we have re-evaluated Figure 2 using a fully independent

236 GSRM simulation dataset, out-of-the-80day-samples. The results (Figure A1) show that the R^2
237 values for Q_1 and Q_2 in this independent dataset closely align with our previous findings in both
238 pattern and magnitude, confirming the good generalization capacity of the trained model.

239 We fully agree that incorporating uncorrelated out-of-sample validation throughout training is
240 essential for quantifying overfitting and identifying the optimal stopping point. Given the time
241 constraints and the model's satisfactory performance in the independent validation described above,
242 we did not retrain the ML model for this study. Nevertheless, the proposed training-to-integration
243 workflow is readily extensible and can be straightforwardly applied in future work with systematic
244 out-of-sample validation.

245 **It is also concerning to have the manuscript state it uses previously-established measures for**
246 **ITCZ strength and width, but then not use the measures described in the cited papers. This type**
247 **of error is easy to miss in review as it is not the reviewer's responsibility to read cited papers.**
248 **It is necessary to use established measures for complex measures like this to avoid p-hacking**
249 **(which I don't think is the case here, the ITCZ representation does look better). Perhaps I am**
250 **not understanding the measures in the cited papers, but this should be explained.**

251 Reply: Thank you for your comment. We acknowledge that the definitions of ITCZ intensity and width
252 in the original manuscript were not sufficiently clear. In the revised version, we adopt more explicit,
253 literature-based definitions with appropriate citations to ensure clarity for you and other readers.
254 Results using the updated definitions are consistent with our original findings and do not affect the
255 study's conclusions. For further details, please refer to our earlier response.