# Response to Referee #3

We would like to thank the three reviewers for their helpful comments and corrections. We propose a revised version of the paper in which we have taken into account their comments and suggestions. We hope that our revisions will address their expectations.

In the following document, the reviewer's comments are shown in **black** and our point-by-point responses are in **blue**.

The main adjustments proposed in the new version of the manuscript are described below:

- The section "Code and data availability" has been modified to add a brief description of the ICOLMDZ model.
- We have reorganized the figures and tables to improve the flow as reviewers suggested. The figure and table numbers have changed from the initial version of the manuscript:
  - Figure 5 has been removed to avoid repetition with Table 5.
  - Table 3, Table 4, Figure 3 and Figure 8 have been moved to the Supplementary Material.
  - Table 6 has been removed.
  - Figures S2 and S4 were moved from the Supplementary Material to the main text.
  - For each tendency, the results obtained during the first learning and those obtained with the second learning were separated to ease the flow of the presentation of the results. The results of the second learning are only shown when they are discussed (in Section 5 "Refined approach: integrating physical knowledge"). The revised figure numbering, for the zonal wind tendency, is as follows: Figure 4 has been split into Figure 3 and Figure 9, and Figure 6 becomes Figure 4 and Figure 10. We did the same for the five other tendencies in the Supplementary Material. We also separated the results from Table 5 into two tables.
- We propose a new colour scheme for vertical profiles according to the Coblis-Color Blindness:

	Initial training	Second training with laplacians
DNN	Green with solid lines	Blue with solid lines
U-Net	Orange with dashed lines	Red with dashed lines

- Section 6 "Discussion" has been merged with Section 7 "Conclusions". The results of this merger is Section 7 "Conclusions and discussion".

- As suggested by the reviewer 3, new results on the emulation of a realistic configuration have been added to the manuscript in Section 6 "Results for a realistic setup". The data used for this realistic configuration are described in Section 2.1 "ICOLMDZ simulation data". The terms "for the aquaplanet setup" have been added to the names of Sections 4 "Initial training performance for the aquaplanet setup" and 5 "Refined approach for the aquaplanet setup: integrating physical knowledge" to avoid confusion.

In the manuscript "Contribution of physical latent knowledge to the emulation of an atmospheric physics model: a study based on the LMDZ Atmospheric General Circulation Model", the authors compare two Neural Network (NN) architectures, Dense Neural Networks (DNNs) and U-Nets, for their ability to emulate tendencies produced by a conventional physics package. They find that U-Nets generally outperform DNNs in capturing the variability of these physical processes. The study demonstrates that feature engineering by integrating physical knowledge, such as including the Laplacians of state variables as additional predictors, can improve the accuracy of both investigated NNs, particularly for the DNN.

I think this research could contribute to the development of more skillfull ML emulators of subgrid physics, but it has significant limitations as it is only trained on data produced by an aquaplanet setup without any topography and the performance of the emulators is only evaluated offline. An evaluation of the online performance, coupled to the aquaplanet ICOLMDZ setup would certainly be more meaningful than the offline evaluation.

We have taken this comment into account and added new results related to the emulation of a realistic configuration in Section 6 "Results for a realistic setup". The data used are presented in Section 2.1 "ICOLMDZ simulation data". The abstract, the introduction and the conclusion have been modified accordingly.

Some of the findings presented in this manuscript appear to be preliminary and would benefit from more thorough evaluation. I will elaborate on this in the following points:

### **Major comments**

1. From my perspective the biggest issue, as mentioned, is that the developed emulators are only evaluated offline. A small stability/performance test of the DNNs/U-Nets coupled online would be very insightful.

We totally agree that online evaluation is crucial. However, we would like to point out that it is out of the scope for this first study because the main objective of this article is to show that adding predictors with physical meaning

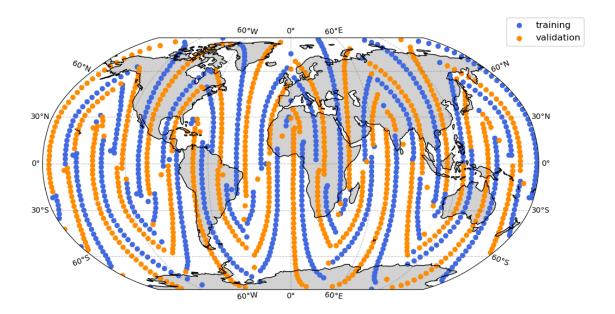
helps to obtain an emulator closer to the physical model. Before studying emulators in an online setup, we need to ensure that they reproduce the behavior of the physics in offline mode. Moreover, online evaluation requires procedures and metrics that differ from offline evaluation and an entire paper would be required to examine this thoroughly. Thus, the online evaluation is presented as a perspective in Section 7 "Conclusions and discussion".

2. Line 160: Why don't you consider 2D outputs such as precipitation. I am sure the physics package of ICOLMDZ also parameterizes some 2D fields.

Our goal was to reproduce the physical tendencies used in the calculations, and thus, we did not consider 2D outputs such as precipitation as mentioned because these variables are diagnostic. This information is now added in Section 2.2 "Input and output variables of emulators" of the revised manuscript.

3. Line 198: Why do you choose one every 20 cells and not a random subset? With this fixed "step size" you potentially overfit to these locations, don't you? What about evaluation on full data set, it is mentioned here, but you never come back to that as far as I saw?

We have chosen one cell every 20 cells because the data is well distributed across the globe, all latitudes are represented, as illustrated in Figure 1 below. This information is now added to the paragraph: "In this way, we obtain data that is well distributed across the globe and we ensure that we have sampled at all latitudes (see Fig. S1 in the Supplementary Material)". The revised version also includes the figure below.



<u>Figure 1</u>. Distribution of grid cells over the globe used for training in blue and validation in orange, for the realistic configuration. The same distribution of grid cells is used for the aquaplanet.

All the results presented in the manuscript come from the test dataset where only temporal sampling is performed.

4. Line 224: You only tried this particular setup (6 layers and 512 neurons in each layer) and no hyperparameter optimization at all?

We tried many different setups for the DNN (and also for the U-Net architecture) and we decided to work with a DNN composed of 6 layers of 512 neurons for each layer since it appeared to be a good trade-off between the computation time and the quality of the results. This information is now clearly added to the paragraph.

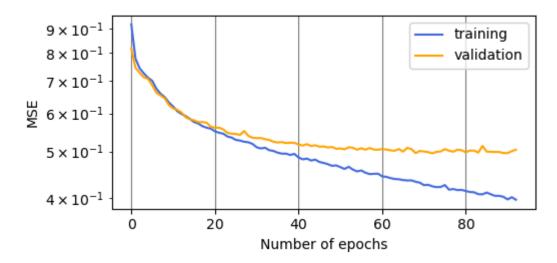
We did not use hyperparameter optimization techniques but we did some manual tuning.

5. Line 483: To judge impact of adding additional predictors adding a plot learning/validation curves to assess convergence would be helpful. Also, best to train multiple models to see if the improvement is "statistically significant"

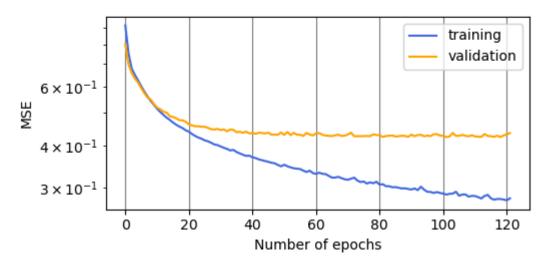
We agree that plots corresponding to the loss function for the learning and for the validation step help to assess convergence of the models. However, we have not added this type of plots to the manuscript because the interpretation is challenging due to the multivariate nature of the emulation problem.

Models do not converge perfectly during training because of the early stopping criterion. This information is now added in Section 3.1 "Neural network architecture": "It should be noted that this stopping criterion may prevent models from converging perfectly, the key point is to avoid overfitting. Thus, the models saved and used are those that achieved the best performance according to the computed loss metric on the validation dataset tracked during training".

For instance, the training and validation loss curves for the DNN without laplacian and for the DNN with laplacians are shown below with Figures 2 and 3, respectively. The early stopping criterion used to prevent overfitting, with a patience of 20, must be taken into account when reading the loss curves. The model selected for the DNN without laplacian achieved a validation loss of 5.0x10<sup>-1</sup>. For the DNN with laplacians, this score is equal to 4.2x10<sup>-1</sup>.



<u>Figure 2</u>. Plot of training and validation loss over epochs for the DNN without laplacian.



<u>Figure 3</u>. Plot of training and validation loss over epochs for the DNN with laplacians.

We have trained each model several times but we have not kept all the corresponding loss values. However, the results are as follows:

- Similar behaviors and performances in terms of loss function for a chosen model have been observed.
- The laplacian-aware emulators achieve a lower loss than the initial emulators without laplacian.
- U-Net models perform better than DNN models in terms of loss minimization.

The following sentence has been added to the manuscript: "Each model has been trained a few times, but as we observed similar behaviors and performance in terms of loss function for a chosen model, we only kept one of each".

6. Line 562: The discussion is very short and not complete, you should add some limitation of you study, e.g. the aquaplanet setup.

We have merged Section 6 "Discussion" and Section 7 "Conclusions" to create a new Section 7 "Conclusions and discussion".

Several discussion points are now provided regarding some elements to understand why the U-Net is more efficient than the DNN, the contribution of the turbulence process for the five other tendencies, and also the execution time estimates.

We are also discussing some perspectives that can be considered as direct limitations of the study, such as the capacity of the emulators to generalize and the online evaluation.

We did not include a paragraph in which we discuss the use of an aquaplanet as a limitation because we have added some results, in Section 6 "Results for a realistic setup", on the emulation of a realistic configuration, even though these are preliminary results.

#### Minor comments

1. Eq. (2): I doubt that the ICOLMDZ uses as simple Euler step for time integration.

Yes, it is an approximation. To clarify, we have decided to remove this equation and the corresponding sentences (lines 134-135).

2. Line 242: Please provide more details how you "reformatted" the input to 2D tensors. You are using 1D convolutions as previous developed ML physics schemes using U-Nets, right?

We have taken this comment into account and we have revised the paragraph accordingly. Here is the information that has been added: "Data is reformatted before being provided to the encoder. For vectors, but also the 3 scalars spread across all vertical levels, we concatenate each atmospheric column in order to obtain a 2D matrix of dimension 79 multiplied by the number of 1D variables. This data format is compatible with the 1D convolution operator that processes each atmospheric column independently".

Yes, we are using 1D convolutions. This information is now clearly added to the paragraph: "We have also used a specific architecture, called a U-Net model (Ronneberger et al., 2015), where we use 1D convolutional layers".

3. Eq. (8-13): I don't think explaining the computations of means, variances, and metrics in such detail is necessary, but that is ultimately your choice.

We have decided to keep these computations in the manuscript.

4. Line 420/Figure 8: Why do you choose a random timestep and not multiple ones and take a median/mean profile? Also, it could be helpful to plot A,B, and C in the same figure for easier comparison.

We thought it would be interesting to have a look at results for some random timesteps before we got an overview of the breakdown. For this reason, we present the results of the breakdown for one timestep in Figures S1, 8 and S3 (now these figures are all provided in the Supplementary Material). The mean vertical profiles, obtained from all vertical profiles over one year, are presented in Figures S2, 9 and S4 (now, these figures are all included in the main text of the manuscript).

If we plot A, B and C in the same figure, it will not be easy to see anything because the scale of the zonal wind tendency differs depending on the point considered: 10<sup>-4</sup> for A, 10<sup>-3</sup> for B and 10<sup>-5</sup> for C, for one time step for instance.

5. Line 423: What physical process is represented by the thermals opposed to convection and turbulence?

To respond to this comment and the following comment, we have replaced "These physical processes include turbulence d\_u\_turb, convection d\_u\_conv, thermals d\_u\_therm, gravity wave drags associated to fronts d\_u\_gwd\_front and those associated to precipitations d\_u\_gwd\_precip" with the next sentence "These physical processes include turbulence d\_u\_turb, boundary-layer convection d\_u\_therm, and convection modelled as static adjustment d\_u\_conv, gravity wave drags due to emission by fronts d\_u\_gwd\_front and those due to emission by convective precipitations d\_u\_gwd\_precip". More details on these processes are given in Hourdin et al. 2020 (3.1. Small-Scale Turbulence, 3.2. Boundary Layer Convection, 3.3. Deep Convection, 3.6. Radiation, 3.8. Gravity Waves).

6. Line 436: "This is also the case for gravity wave drags from precipitations" - Do you mean convectively coupled waves, or what should "gravity waves from precipitation" mean?

Please, see our response to the previous comment.

7. Line 470: Is delta z really fixed and if so, what is this fixed distance?

 $\delta z$  is fixed and is equal to 1. It corresponds to a difference between levels (and not to a difference in altitude). We have modified the text and here is the revised version we propose: "In the denominator, delta z corresponds to the difference between vertical levels. We consider  $\delta z = 1$ . This distance is fixed as the step between vertical levels is constant".

8. Line 613: You should mention the CPU inference time of your models as well and mention that even if you are able to use the GPU for inference during coupling, you have an overhead from copying the data from CPU to GPU as long as ICOLMDZ is not running on the GPU as well.

We have modified the paragraph concerning the execution time estimates by adding the following sentences: "However, we would like to emphasize that estimating the gain when using an emulator instead of a physical model is much more complex. This is why these estimates should be considered as a preliminary test. It remains to be seen whether this gain would remain when the emulators are coupled with the dynamical core of the model, i.e., in an online setup. For instance, there is generally a substantial overhead when exchanging data from the dynamical core run on CPUs to the physics emulators run on GPUs. One way to alleviate this issue would also be to run the dynamical core on GPUs, which is already possible with the latest development of DYNAMICO".

 Line 631: Watt-Meyer et al., 2024 did not emulate a physics package but derived the subgrid physics by coarse graining from high-resolution simulations

Thanks for pointing this out, we apologize for the confusion. We have removed the quote in question.

## Typos/Grammar

- Line 19: The second "numerical" is redundant Fixed.
- 2. Line 201-203: "These numbers of samples are multiplied by the length of the input vector X' and the output vector Y', for all the three subsets, which correspond to variables aggregated over the vertical dimension" What do you mean by multiplied by length of ..., to get the number overall size of the input data?

Exactly, we added this sentence to illustrate the size of the data. However, in the new version, we have decided to remove this sentence to improve clarity.

3. Line 229: "It is trained with the learning rate scheduler of 10^{-3} ..." What exact scheduler did you use? I assume the factor 10^{-3} indicates some decay or did you mean learning rate of 10^{-3}?

Our apologies, it is a mistake. The results presented in this manuscript come from models with a fixed learning rate, we did not use a scheduler here. Therefore, we have removed the term "scheduler".

4. Line 230: I wouldn't say that "a U-Net is a CNN architecture" but that it uses convolutional layers

The initial sentence has been changed by "We have also used a specific architecture, called a U-Net model (Ronneberger et al., 2015), where we use 1D convolutional layers".

5. Line 246: "Once the models had been built, we trained them." Not necessary to state this.

Noted, we have removed this sentence.

6. Line 257: "In this article" - You mean section I think

We have replaced "In this article" with "In the following sections".

7. Line 284: "We will sometimes compute the RMSE on specific subsets of the test dataset to exhibit performances that vary for instance across latitudes or vertical levels." - "exhibit" does not fit here

Corrected, we have replaced "exhibit" with "assess".

8. Line 333: "..., we examine them study by latitude belts, ..."

This has been corrected, we have removed the term "study".

9. Line 416: "We decided to carry out this study on three grid cells..." - please use columns instead of grid cells here

This has been corrected as suggested by the reviewer.

10. Line 429-431: please rewrite the sentence

The initial sentence was "Turbulence movements (see Figure S1a, 8a and S3a) are present whatever the point studied near the surface but also at higher altitudes, between vertical levels 20 and 40, for points located at high latitudes and in the subtropics". We have changed this sentence and here is the new proposal: "Whatever the point considered, turbulent movements (see Figure S3b, S4b and S5b) can be observed near the surface. These movements also occur at higher altitudes, more precisely between vertical levels 20 and 40, for the points located at high latitudes and in the subtropics".

11. Line 535: reformulate "with these new learning"

Done, it has been replaced by "with the laplacian-aware emulators".

12. The word "Indeed" is overused throughout the manuscript.

We have revised the text to remove several occurrences of "indeed" (9 terms in total now, compared to 18 initially).

## **Reference**

Hourdin, F., Rio, C., Grandpeix, J.-Y., Madeleine, J.-B., Cheruy, F., Rochetin, N., et al. (2020). LMDZ6A: The atmospheric component of the IPSL climate model with improved and better tuned physics. Journal of Advances in Modeling Earth Systems, 12, e2019MS001892. <a href="https://doi.org/10.1029/2019MS001892">https://doi.org/10.1029/2019MS001892</a>