

## # General comments

RC: This is a very interesting study that tackles the difficult task of simulating N<sub>2</sub>O emissions using LSTM models across European grasslands. I think that the dataset compilation is a valuable contribution to the community. To help this manuscript reach its full potential, I have identified a few key areas where the approach should be strengthened regarding data aggregation, validation strategy and data transformation. In addition, the Discussion would benefit from a stronger connection between findings in this study with existing literature. I believe this extra effort is necessary to prove the robustness of the presented approach and will help this paper to make a significant contribution to this field. For this reason, I recommend major revisions.

AC: Thank you for your detailed and constructive review. We appreciate your feedback and implement your suggestions. Specifically, regarding the validation strategy (data split): We improved the train/ test split to ensure there is no doubt of data leakage. Regarding data transformation: we replied to your comments about log-transformation and loss function to clarify the motivation and correct implementation. Regarding model performance analysis, we added new metrics regarding the model performance and made sure we report specifically. We refined the discussion by strengthening the context of our findings with the results of existing publications. We hope that these improvements and revisions will meet the high standards of Biogeosciences and provide an interesting, publication-ready manuscript.

## # Major comments

### ## Aggregation to daily values

L70-97: The manuscript states that sub-daily flux measurements (from both automatic chambers and eddy covariance) were "aggregated into daily values". However, the specific method of aggregation and the quality control criteria for this aggregation are not defined here. In addition, I am missing info about which data were aggregated (gap-filled or non-gap-filled). Especially for the eddy covariance (EC) sites Chamau and Neustift this would be crucial information for reproducibility, since half-hourly EC data are frequently rejected during quality control which often leaves days with very few valid data points. Then, calculating a daily value from only a few sparse measurements (which may be biased towards specific times of day) can introduce significant errors. It is therefore currently unclear if the daily training targets represent true daily averages, or if they are biased artifacts of incomplete data.

I therefore suggest that the authors please specify the aggregation method and give more detailed info about the used time series (gap-filled or non-gap-filled) that were aggregated. In case non-gap-filled time series were used, the minimum data availability threshold required to accept a day as valid should be given. If days with low coverage were included, please discuss the potential bias this introduces to the model. If gap-filled time series were used for the aggregation, the used gap-filling models should be given.

AC: For N<sub>2</sub>O, due to its sporadic nature gap-filling using temporal interpolation or look-up tables can be misleading, since preceding fertilizer events are usually not considered and the period before a fertilizer event can be erroneously filled with too high values. We wanted to avoid comparing ML model results to modelled data, which seems to be the case if larger gaps in the sub-daily timeseries were filled with a more sophisticated empirical model, we prefer to compare the ML model only with direct measurements. Since the investigated N<sub>2</sub>O fluxes showed no diurnal cycle we consider this approach justified and implemented it for the Tereno data which we directly processed. We added the data QA/QC was performed by the experimentalists publishing the original data. Due to differences in how we received the data from experimentalists, the approach of aggregating sub-daily to daily data differed across sites.

1. For Chamau, Switzerland: EC post-processing was done using EddyPro, and only fluxes labelled with quality flags 0 and 1 were included (See Fuchs et al. 2018; For details section on QA/QC and post-processing). Days with at least 30% measurements were used. Diurnal cycles were absent, thus an averaging to daily values without previous gap-filling was justified.
2. For TERENO, Germany: Experimentalists went through data collection: The robot places the chamber every 15 minutes on a lysimeter. However, if the wind is too strong, the robot cannot lower the chamber correctly, so that it may go into park position at high wind speeds. There is a higher absolute measurement frequency in Graswang compared to Fendt.

We received the data which included up to 6 measurements per day and we included all days with at least 2 out of 6 possible measurements on that day.

3. For Neustift, Austria: EC post-processing and gap-filling of sub-daily data was performed by Hörtnagl, L., & Wohlfahrt, G. (2014): „Gaps less than or equal to 2 h were filled by linear interpolation. For the filling of larger gaps a lookup table was generated, using flux data in a time window of 14 days around the missing flux value and T<sub>soil</sub> bin widths of 1 °C.“ We used this data that was gap-filled and then aggregated to daily values.

discussion section 4.5 Implications and limitations : L398-L400, Additionally, the number of available sub-daily flux measurements varied across days in the TERENO and Chamau non-gap-filled datasets due to gaps resulting from quality control procedures and instrumental failures (Kiese et al., 2018; Feigenwinter et al., 2023). Such uneven data availability can introduce bias into flux estimation modelling (Lucas-Moffat et al., 2022; Teng et al., 2020).

### ## Input variables

L103-114 and Table 4: The model looking back at drivers of the past 5 days is likely too short to capture the causal drivers of emission peaks, because the soil is potentially primed by an event long before the emissions actually happen (legacy effect). For example, if

fertilization happened 6 days ago then the model cannot see it and it assumes e.g. "rain and no fertilizer" in its time window. Therefore, a longer sequence would be more appropriate (e.g. 30 days or longer). One architecture was tested using 10 days, but this might also be too short. Maybe the authors have already tested longer windows (10+) and can extend Table 4 with this additional information. It is still possible that the shortest window works best, but it is important to see this in comparison to other tests.

Thank you for the valid point. Initially, we started with sequence length 5, then we checked if longer sequences had an effect on model performance. Several points led to the decision to remain with shorter sequences in the publication:

(1) Changing the sequence length to 10 (and hyperparameter tuning) or 20 didn't bring significant improvements to the model so we didn't investigate longer sequences. Please look at the Table 1 below showing how the model acts on samples of test set when we include a longer sequence of length 20:

```

=== LOG SCALE ===
                R2 Spearman Spearman_CI 95% RMSE MAE Bias
Dataset
Train          0.536    0.587    [0.58, 0.60] 0.733 0.604 0.063
Test           0.380    0.570    [0.55, 0.59] 0.693 0.542 0.352
Independent test -0.328    0.361    [0.34, 0.38] 0.715 0.607 0.401

=== REAL SCALE ===
                R2 Spearman Spearman_CI 95% RMSE MAE Bias
Dataset
Train          0.346    0.587    [0.58, 0.60] 43.570 12.469 -6.292
Test           0.417    0.570    [0.55, 0.59] 26.806  8.501  1.740
Independent test -0.154    0.361    [0.34, 0.38] 11.677  4.925  1.947

```

Table 1. Statistical metrics on a model trained on sequence length 20.

(2) Longer input sequences (length **X**) increase the amount of data required for the model to function effectively, which makes evaluation at a new site more challenging. This is because the LSTM requires **X consecutive days of input** to estimate the **N<sub>2</sub>O flux for a given day**. As the required sequence length increases, the amount of usable data decreases due to gaps in input variables such as soil moisture and temperature. Consequently, we chose to keep the sequence length—and thus the model's data requirements—as small as possible to ensure robust performance. The potential benefits of longer sequences (e.g., **30 days or more**) will be investigated in future work using an expanded dataset.

we also added a statement in line L401:L403 in the Implications and limitations section to ensure using a longer sequence is known as a potential future work: „Firstly, expanding the dataset to encompass more geographically and climatically diverse sites, along with a wider range of management practices and legacy effects on N<sub>2</sub>O emissions, would allow for the development of more generalized N<sub>2</sub>O emission models for grasslands across different agroecosystems“.

### ## Dataset splitting

L120-130, Data leakage: The manuscript describes a data splitting strategy where parallel treatments (e.g., intensive vs. extensive management) from the same site and year are separated into training and testing sets. However, this approach introduces significant data leakage because the parallel treatments experience identical daily meteorological drivers. Basically, the model is exposed to the (specific) weather sequences of the test period during training. I would assume that performance metrics on the test set likely reflect the ability of the model to memorize these weather patterns rather than the true ability of the model to generalize. A test for generalization would require the test data to be independent in time (unseen years) or space (unseen sites).

I therefore suggest that the authors restructure their validation strategy. This could be done by withholding entire years across all treatments or by withholding entire sites (as was done with Rottenbuch) from the training set.

RC: Thank you for your feedback on our validation strategy and train-test split. We restructured the concept and addressed your concern about data leakage. In the updated version we have a clear definition of a cross-site Test Set for independence in space (Rottenbuch, intensive as well as extensive lysimeters) and a Temporal Test Set for independence in time. We specify this in the paper in L126-L140.

L125/L130, Train/test splitting:

In Section 2.3, the reported train/test split (13000 training vs. 15000 testing) is highly unconventional for machine learning and risks underfitting. It is unusual for a test set to be larger than the training set, I think I have not seen this in other studies yet. Unless there is a very specific requirement for these unusual data sizes that I am not aware of (and a proper justification), this needs to be addressed in an updated version of the manuscript.

RC: The reason for a huge test set is considering a cross-site test that has many years of data in total 8690 data rows. Our new train-test split results now in a more conventional split (14500 for training vs. 8690 for cross-site testing and 3160 for temporal testing).

### ## Log-transformation of N2O data and loss function

L140-157, L178-189: The text describes that daily N2O fluxes were log-transformed, meaning that the model outputs predictions on a log scale which are then back-transformed. This transformation seems to be somewhat in conflict with L180 that emphasizes the "importance of catching high peaks". To me it seems counterintuitive to first compress extreme values/peaks (log transformation), but then applying a custom exponential loss function that basically attempts to re-weight the peaks. By compressing the target variable, it incentivizes the model to focus on fitting the "background" noise of low magnitude emissions rather than learning the rare, high magnitude events that drive the N2O budget. Therefore, the ability of the model to predict timing and magnitude of emission peaks is limited. Regarding the loss function, as pointed out the penalty grows exponentially as the error increases. If used in a dataset with rare high peaks this makes the model highly unstable during training.

I am sure that there were good reasons why the authors choose this approach, and I therefore suggest that the authors add their reasoning to the text. If the authors want to reconsider their approach and aim for optimizing the predictive power for peak events, I recommend to train the model on the linear N2O data in combination with a loss function that is robust to outliers (to improve stability).

AC: We chose this approach due to the following reasons:

(1) Applying the Log-transformation on highly right-skewed data is a crucial step before applying a model: The log transformation reduces skewness and stabilizes variance, resulting in better-behaved gradients, faster convergence, and more stable and effective training (West 2021; Bishop, 2006, Ch. 5; Hastie et al., 2009, Sec. 2.4; Box & Cox, 1964). In highly skewed data, rare extreme values dominate the gradient in gradient-based models such as neural networks, which is known to impair optimization and slow convergence in deep neural networks (Goodfellow et al., 2016, Ch. 8). Variance-stabilizing transformations such as logarithmic scaling are commonly used in recurrent neural forecasting models to address strong skewness and scale heterogeneity (Salinas et al., 2020; Smyl 2020). Without log-transformation, the extreme values create high gradients, dominate the training, and cause the model even not learn the relationship in the data but focus on the rare 5% of the extreme events only. Our goal is not to only predict large extremes, but learn the N2O emission pattern, shown in both low and high fluxes. We briefly mentioned the importance of log-transform in L150:L153

(2) Loss function: The log-transform helps training stability which is a must in this problem as discussed, but we still needed a proper loss function for the specific problem. Our goal is to have a model optimized to learn both high and low emissions, rather than only extremely high values. A weighted loss function ensures a trade-off between the target transformation for better convergence, and high performance of the model on high emissions. In addition, the exponential loss function is executed on log-transformed data, making it not only penalize large errors on high values, but also on low fluxes, which ensures the model learns both.

(3) Training stability: Exponential weighting functions may introduce numerical instability if poorly constrained. But, to ensure stability, model inputs were scaled, the peak-weight parameter was conservatively bounded (1-1.5), and training convergence was carefully monitored. No gradient divergence or optimization instability was observed. The exponential weighting therefore functioned as a controlled mechanism to reduce bias toward the mean and improve the representation of extreme events, which constitute a primary objective of the study.

(4) The loss function in our study was designed to emphasize rare high-magnitude flux events, which are systematically underweighted by conventional losses such as MSE or MAE in highly right-skewed datasets. MSE, MAE, Huber loss, weighted Huber and weighted MSE were also tested, which all underestimate high peaks and/or overestimated low fluxes with a flat line threshold.

(5) The use of standardization and normalization in combination with different loss functions such as MAE, weighted MSE, Huber loss as well as the custom loss was tested. Indeed huge amount of experiments using approaches which theoretically seem to be a proper fit for this problem were tested. Finally, the investigations showed that log-transform was the best approach, also underpinned by all the aspects we discussed above in 1).

An example of another study which used log-transformation and a weighted loss function on highly skewed data was Ghosh et al., 2025: "Short-term precipitation prediction in India using ConvLSTM", which used weighted mean squared error loss function.

### ## LSTM

L163-175: In its current state, Section 2.4.1 lacks technical specifications that would enable others to reproduce the model. The parts in the text that mention "a set of LSTM layers" and "several dense layers" are missing details, e.g. the specific layer dimensions (units per layer) used in the final model. Similarly, the "attention mechanism" should be explained, at least briefly, even if more details are given in the provided reference. Also, attention can be applied in different places in the pipeline and there are multiple ways for its implementation. These details, among others, are currently missing, but required for reproducibility. I also suggest to add some of these details to Fig. 2.

AC: In order to increase clarity, we now put the technical specifications of the LSTM also in the method section 2.4.1. In the previous version we had put it only in the result section: section 2.4.1 LSTM model architecture

. To ensure reproducibility and full documentation, we put the code to Github, and the link could be found in the Data and code availability section. This code includes the loss function, and the LSTM model and specific values of hyperparameters which are finally selected for the model.

### ## Hyperparameter tuning

L190-200: I appreciate showing the ranges used for hyperparameter tuning, but this section also needs to show the optimized values that the authors decided to use in the final model. In addition, in this section I have another concern regarding validation. From the description it is not clear which data were used for validation, so I assume the test set was used instead of a validation set. Section 2.3 only mentions a test set. A typical split for train/validation/test split would be e.g. 60/20/20. To my knowledge, optimization algorithms like Hyperband require a validation metric to compare the different configurations. If the test set was used for optimization, this would be a "training on the test set" error (it would mean the model has seen the "answers" to the tests it is taking), a form of data leakage. If the training set was used for optimization, the hyperparameters are likely overfitted to the training data. And finally, the "manual adjustments" are in my view difficult to justify after the optimization procedure, because it means that the model was optimized, but optimal parameters were then not applied. Generally, manual tuning based on observing model output introduces researcher bias and implies that the model was tweaked until it \*looked\* right on the evaluation data. Regarding the ranges for tuning: it is important that the authors justify the specific ranges chosen and describe e.g. why the search space for neuron count was so constrained (it is very narrow and specific).

AC: (1) After the train/test split which is mentioned in section 2.3, the test set is completely hold out of the training, the training set is partitioned further to train/validation with 80/20 ratio. So, no info from the cross-site or temporal test gets into training According to your feedback, we add a short and clear explanation about it in L 130-132: „During the model training process, the training dataset partitions into training and validation parts following an 80:20 split as a requirement for machine learning model training. The training partition is used to fit the model, while the validation partition facilitates hyperparameter optimization—such as learning rate decay—and helps monitoring model's progress during training.“

(2) Regarding manual hyperparameter tuning: The hyperparameter search space and manual adjustments are required to ensure reproducibility. Importantly, adjustments were performed exclusively using the training and validation datasets to avoid information leakage, and final model performance was assessed on an independent test set. HPOs provide a great opportunity toward finding a good model, but none of them are perfect (Raiaan et al., 2024). Manual adjustments, guided by domain expertise and iterative experimentation, can further refine model performance by incorporating nuanced understanding of the problem context and model behavior. Manual adjustments included small changes in dropout, learning rate, loss function coefficient and batch size. For instance, we changed dropout from 0.05 to 0.03 or batch size from 120 to 100 to see the potential improvements. The value of 0.03 was not in the HPO search space, but as 0.05 was in one of the best configurations returned by HPO, so we tested the values around it like 0.03 and 0.07. The number of neurons and layer regularization were retained as determined by Hyperband. Without these adjustments, the model tended to generate predictions around center of distribution that yielded lower statistical error metrics but failed to accurately represent abrupt increases or decreases in emissions, which constituted an important objective of N2O emission estimation. Ultimately, model development required careful trade-offs among the loss function design and weighting, batch size, learning rate, and dropout. Compared with typical modeling tasks, this process demanded substantially more experimentation, including the evaluation of alternative data-scaling strategies.

(3) Regarding narrow neurons size: the ranges of hyperparameters to go through hyperparameter tuning is a decision based on expertise. We also tested a higher number of neurons in few experiments, but that didn't increase the performance and the search space for the number of neurons was then narrowed for HPO algorithm and higher number of neurons were excluded from the rest of experiments. For other hyperparameters, the ranges are best practice values which are always tested for LSTMS.

### ## Feature importance

L202-209, Fig. 6, Fig. 7: PFI is used to assess the importance of drivers at individual time steps. However, this is methodologically flawed for highly autocorrelated time series such as soil moisture since the value at time  $t$  is nearly identical to  $t-1$ . So by shuffling only the value at time  $t$  the model can recover the information from the unshuffled values at  $t-1$  (or  $t-2$  ...) stored in the memory of the LSTM. Consequently, the importance score of the respective feature is reduced and continuous drivers might be undervalued in comparison to e.g. fertilization. Therefore, I am not sure Fig. 7 can work because the lags are not independent. For Fig. 6, it is not immediately clear to me why the lagged variables are missing. I assume that maybe all time steps of a single variable were shuffled together, and then an overall importance is given (sort of a "grouped" PFI based on grouped shuffling)? After getting comments about the PFI for time lags, and considering different opinions on this topic between data scientists which doubts whether this is a proper approach or might cause misunderstanding about model's results, we decided to remove Fig. 7 and solely rely on Fig. 6 as a standard way which the PFI is calculated in other scientific publications. We believe that such a method needs very wide discussion about all the aspects and the implementation method, which its out of the scope of this paper.

In Figure 6, shows if we initially shuffle all the values of a variable, regardless of lags and then execute the model on it, so what you mentioned is correct, exactly the overall importance of that variable is considered.

### ## Model performance

L230 Table 4: The relative RMSE of 270% is very high, and one of the reasons could be that the log transformation distorts the optimization. The log transformation minimizes the (logarithmic) error, but after the results are back-transformed to the linear scale the bias on low-flux days becomes disproportionately large (relative to the mean). Also, I assume the errors on peak events increase dramatically. The authors should compare this performance against e.g. the predicted mean of the training set for every day and then check whether the LSTM can outperform this naive baseline.--> we added  $R^2$  for this case.

In light of the high RMSE, it is important that the authors discuss why the error is so high, especially if it is driven by a few very high outlier values, or if the model is systematically biased. A scatter plot of predicted vs. observed fluxes on the linear scale (not log-log) can help here to diagnose this. Currently, I cannot see how the model is "capable to transfer results across sites".

AC: The RRMSE for N<sub>2</sub>O predictions usually largely depends on the dataset. Since the RRMSE divides by the site mean, a low divisor creates a high RRMSE. Due to the sporadic nature of N<sub>2</sub>O which creates a highly skewed distribution of N<sub>2</sub>O values, the reported RMSEs and also higher values than usually reported for predictions of other greenhouse gases like CO<sub>2</sub> (and CH<sub>4</sub>). Only few extreme errors combined with very low mean of a highly right skewed dataset highly increase RRMSE and R<sup>2</sup> scores, thus the model's performance can not be judged only by these metrics. Therefore, also by considering your later comment on the fact that relative RMSE is probably not a good metric for this dataset, and the comments of other reviewers in requiring more metrics, we improve our model performance analysis by reporting several metrics which could help better analysing the results. Please see Table 5 in section 3.2 model performance for detailed result analysis on the train set, temporal hold out test set and cross-site test set.

Regarding log-transformation and specifically its importance on highly right-skewed datasets to stabilize training, we refer to the answer earlier in this document.

#### # Specific comments

L10/L61: The site Neustift is described as "North Austria". Please correct this to "Western Austria". [Corrected](#).

L25: Change "greenhouse gas" to "anthropogenic greenhouse gas". [Done](#).

L30: The sentence "The Pearson correlation..." presents specific results in the introduction ("between -0.11 to 0.16 for our dataset"), which should not be done. I suggest the authors rephrase this sentence and remove these specific calculation outcomes, they have their place in the Results section.

[We removed this from the introduction and added to the results, L223-L225](#)

L33: Change "long periods" to "long time periods"

[Done!](#)

L56: I would not go so far to mention "general applicability". One independent grassland is not enough for this claim. [Changed to „regional applicability“](#)

L62: "no grazing": Some of these sites have grazing. If the authors are referring to specific time periods used in their study, please add this info.

[Thanks a lot for mentioning this- We'll correct this in the paper. Chamau had some days grazing in 2014-few days in 2015 and some days in 2016. We remove phrase „no grazing“ from L62. Also in section 4.5, L434-435 it mentioned that: „For instance, different grazing practices could affect the N<sub>2</sub>O emission patterns, which were not considered in this study due to limited data \(only in Chamau there were 3 years including grazing practices and other sites in this study were non-grazed\).“](#)

L65: I suggest to give the FLUXNET identifiers for all sites, which in case of Chamau is "CH-Cha" (instead of "CH-CHA", note the lowercase letters) and in case of Neustift "AT-Neu" (instead of "AT-NEU"). [Done!](#)

L65 table caption, and other places throughout the manuscript: "Hörtnagel et al., 2014": the correct name is "Hörtnagel" [thank you, corrected!](#)

L65 table header: I think there is a typo in the header and "(kg h-1 y-1)" should be "(kg ha-1 y-1)". I would interpret "h-1" as "per hour". [You are definitely right, thanks for noticing this, corrected!](#)

L65 Table 1: Please give full units for "Annual N<sub>2</sub>O", i.e., is this "kg N ha-1 yr-1"? The header indicates that it is "kg N<sub>2</sub>O ha-1 yr-1", because the header mentions specifically N<sub>2</sub>O, but from the magnitude of the values I would assume N<sub>2</sub>O-N. [Its is indeed N<sub>2</sub>O-N, corrected!](#)

L65/L100 What is "Dataset size"? I am also not sure how this dataset size matches up with the dataset sizes given in Table 2, e.g., Graswang has dataset size 8708 in Table 1, but numbers in Table 2 add up to 1670. Please check and add missing info.

[The dataset size refers to the number of daily values used, thanks for spotting the inconsistency here. As the size of data after pre-processing is given in table 2, we remove repetitive info from table 1.](#)

[It was a mistake in table 1. Initially we calculated all the data received for each grassland before pre-processing to report available data. In Graswang the records of soil moisture and temperature and climate are there for 2014-2020, while N<sub>2</sub>O is only given in 2020 \(as mentioned in Table 2\). This was mistakenly not considered in the initial counting.](#)

L71: For the three German sites, please add "respectively" after giving the altitudes.[Done!](#)

L71: "asl": I suggest changing to "a.s.l.", because this journal generally uses periods for text abbreviations. [Done!](#)

L81/82: The given references point to subpanels a and c in Fig. A1, but this figure does not show any subpanel letters. [Solved!](#)

L85: "is part of the FLUXNET initiative": I suggest to change to "is part of Swiss FluxNet". [Done!](#)

L88: "by separating two parcels at north and south of the site" change to "by establishing separate parcels in the northern and southern parts of the site", if I understand this correctly.[Done!](#)

L89: "in south" should be "in the south" Done!

L100 Table 2: Great to see (and have available) the different soil properties for the different locations. Would it make sense to include the soil properties as static drivers in the model? It definitely could make sense, mainly for an expanded dataset and future works. For this dataset, due to high correlation between some of the soil properties only pH and BD are used. For instance, SOC and BD are highly correlated in the current dataset. Also, clay content and BD showed a strong negative correlation.

L116: "among German grassland" should be "across German grasslands". Done!

L117: "capable to transfer" should be "capable of transferring" Done!

L125: "13000 data sequences": to me it is not immediately clear what "data sequence" means. I assume one data sequence is what is shown in Eq. 1, i.e., a single input block that is fed into the LSTM (correct?). I suggest to add more info here since the term "data sequence" is mentioned for the first time.

That correct, we meant the available number of sequences that the model looks at. Because, due of gaps of input parameters in some days, the time series were not continous. But, you are correct, this makes a confusion. We use the number of data rows after removing gaps instead to avoid any confusion. Change to „data rows“

L130: "sequence" should be "sequences". Done!

L161: "should be considered" is prescriptive language and OK in other places. In the Methods section I suggest to focus on reporting what "was" done to clearly describe the specific actions taken in this study, not what "should" be done. Done! We changed it to „were considered“

L213: Please add this reference for Matplotlib: <https://doi.org/10.1109/MCSE.2007.55> Done!

L235: Given the low flux values (in combination with sometimes emission peaks), I do not think that the \*relative\* RMSE is a good indicator

We agree. We changed this to other indicators, please look into section 3.2 Table 5.

L235: From the figures I take that flux values are expressed using N<sub>2</sub>O-N, but in the text there are several places where this is not clear. I suggest to add a sentence in the Methods that fluxes are always expressed using N<sub>2</sub>O-N. In this regard, the header in Table 1 is slightly misleading because it mentions "Annual N<sub>2</sub>O" with incomplete units, but the numbers seem to give cumulative N<sub>2</sub>O-N.

Corrected the table 1 as well as anywhere a flux amount was mentioned. Also added this to L65: The amount of fluxes is always expressed as N<sub>2</sub>O-N amount.

L236: Please explain "total bias"

I think you mean L246. we meant on the whole test set, not individual sites. We now report metrics differently for each part of test set as in table 5 and section 3.2

L260: Is there a reason why the time lagged variables are shown in Fig. 7, but not in Fig. 6? Maybe I have missed this somewhere.

Two different approaches of studying feature importance were taken to show feature importance from different perspectives. No other reason behind.

L275 Fig. 4: Typo "occurance" --> "occurrence"

Done!

L290-296: This could be shortened, most was already mentioned before.

We shortened this section, please look into lines L316:L325

L297-306: This currently reads more like a literature review than a discussion. The limitations of previous studies are listed, but the direct connection to findings in this study are missing. I suggest the authors revise this section to synthesize these citations with their results. For example, how their study addresses (or contrasts with) methodological gaps they identified in prior work (e.g., lack of independent test sites).

We improved this section mentioning the methodological gaps that we are filling, please look into lines L316:L325

L308-326: Although this is part of the Discussion, there are no references to existing literature and many of the sentences belong to the Results section. Please revise this section and contextualize findings by citing previous work that supports or contradicts results from this study. For example, is the observed moisture dependency consistent with established theory?

In the revised version section 4.2

L330: "Soil temperature, ..." This sentence seems to be missing words, please fix. [Done!](#)

L331: I am critical of the claim that maximum air temperature can substitute for soil temperature with "minimal impact". Air and soil temperatures often decouple due to factors like snow cover (insulation) and vegetation shading. And of course, there are different soil depths and different thermal lags. In case a sensitivity analysis was performed to demonstrate that the performance drop is indeed negligible, it could be mentioned solely in the context of this study rather than implying it is a viable strategy for sites globally. That's totally a valid point, but in our case we actually had a high correlation between soil temperature and tmax, and we tested creation of the model with either of them and didn't see any significant performance drop when using tmax instead of soil temperature. Therefore, in the paper we also emphasized that in case of a high correlation they could be replaced: „but in case of high Pearson correlation with maximum air temperature (0.86 in our dataset)“, the two variables could be replaced.

L334-337: Context to literature is missing [Done!](#)

L338-345: Here, some interesting points are raised from the literature (pH, BD), but they are not discussed in relation to findings in this study. Instead, the final two sentences jump back to soil moisture but without giving references.

we expanded this section with adding a conclusion sentence and more references please see L381:L387 . In general for each input variable one short paragraph is discussed.

L346-369: I suggest to limit this section to sites that were also used in this study. A comparison with sites from other locations is difficult, especially for N2O that is characterized by high spatial variability.

considering the fact that we got critics about this section from all reviewers, who also questioned if we even need this section, we remove it from this publication!

References:

Goodfellow, I., Bengio, Y., & Courville, A. (2016). *Deep learning*. MIT Press.

Ghosh, T., Anand, S., Nannewar, R. G., & Nagaraj, N. (2025). Deep Learning for Short-Term Precipitation Prediction in Four Major Indian Cities: A ConvLSTM Approach with Explainable AI. *arXiv preprint arXiv:2511.11152*.

Hastie, T., Tibshirani, R., & Friedman, J. (2009). *The elements of statistical learning: Data mining, inference, and prediction* (2nd ed.). Springer.

Box, G. E. P., & Cox, D. R. (1964). An analysis of transformations. *Journal of the Royal Statistical Society: Series B*, **26**(2), 211–243.

Bishop, C. M. (2006). *Pattern recognition and machine learning*. Springer.

Salinas, D., Flunkert, V., Gasthaus, J., & Januschowski, T. (2020). DeepAR: Probabilistic forecasting with autoregressive recurrent networks. *International Journal of Forecasting*, **36**(3), 1181–1191. <https://doi.org/10.1016/j.ijforecast.2019.07.001>

Smyl, S. (2020). A hybrid method of exponential smoothing and recurrent neural networks for time series forecasting. *International Journal of Forecasting*, **36**(1), 75–85. <https://doi.org/10.1016/j.ijforecast.2019.03.017>

Raiaan, M. A. K., Sakib, S., Fahad, N. M., Mamun, A. A., Rahman, M. A., Shatabda, S., & Mukta, M. S. H. (2024). A systematic review of hyperparameter optimization techniques in Convolutional Neural Networks. *Decision Analytics Journal*, **11**, 100470. <https://doi.org/10.1016/j.dajour.2024.100470>

West, R.M. (2021) Best Practice in Statistics: The Use of Log Transformation. *Annals of Clinical Biochemistry: International Journal of Laboratory Medicine*, **59**, 162-165. DOI: 10.1177/00045632211050531.