

## Response to Reviewer 1

*Overall, this manuscript develops a deep-learning model (TCG-Net) for tropical cyclone genesis (TCG) prediction using hourly reanalysis data. The topic has potential practical value, and the paper reflects substantial effort in data preparation and model design. However, the main limitations are that the evaluation framework and baseline comparisons are not sufficiently comprehensive, making the benefits and trade-offs of the proposed approach relative to traditional methods unclear. In addition, the treatment of uncertainty in the “genesis time/location” definition from best-track data is insufficient and may be inconsistent with an hourly prediction setting. The manuscript also lacks case-level analyses and more direct XAI diagnostics to support physical consistency and conclusions about variable contributions. Given these issues, I recommend that the authors strengthen the work by adding more systematic evaluations of metrics, baseline comparisons, sensitivity to genesis definition, case studies, and explainability analyses. Detailed comments are provided below.*

**Authors’ response:** We thank the Reviewer for the positive feedback and constructive suggestions. In this revision, we have carefully addressed your comments, which are detailed in our point-by-point responses below. We hope that the revisions adequately resolve your concerns and further improve the clarity and contribution of our work.

### *1.1. Evaluation method*

*The authors use three metrics—precision, recall, and F1—to evaluate their models. Precision and recall can be strongly affected by class imbalance, and F1, as a function of precision and recall, inherits similar limitations. The authors should include additional evaluation metrics, such as the ROC curve (and AUC) and the precision–recall (PR) curve (and PR-AUC), which assess performance across thresholds and are more informative under imbalanced settings. In addition, it would be useful to test the model with randomly sampled negative examples to examine whether TCG-Net remains robust and whether its performance depends on the specific negative-sampling strategy.*

**Authors’ response:** We appreciate Reviewer 1 for highlighting the limitations of threshold-dependent metrics under class imbalance, for which we fully agree. Following your comments, we have extended the evaluation framework in this revision to include two new threshold-independent metrics, i.e., ROC–AUC and PR–AUC, in all of our analyses. As shown in our revised figures 3, 4, 5, 8, and 9, the results show that ROC–AUC remains consistently high across forecast lead times in both labeling methods (approximately 0.76–0.77 in the Past Domain and 0.79–0.87 in the Dynamic Domain method), indicating stable discriminative ability. Of note, PR–AUC, which is more sensitive to class imbalance, demonstrates consistent performance and confirms that improvements in our ResNet-18 design are not artifacts of a specific threshold choice. For the Past Domain method, PR–AUC increases progressively with lead time (from 0.06 at 6 hours to 0.27 at 48 hours), while it remains lower but consistent with the higher prediction difficulty for the Dynamic Domain method. These additional metrics thus reiterate that our ResNet-18 model maintains strong ranking capability and reliable positive-event detection under imbalance. We have revised the manuscript accordingly with all of these additional ROC and PR curves and analyses.

*1.2. In addition, it would be useful to test the model with randomly sampled negative examples to examine whether TCG-Net remains robust and whether its performance depends on the specific negative-sampling strategy.*

**Authors’ response:** Your point is well taken. In response to your suggestion of evaluating robustness against the negative-sampling strategy, we have revised our analyses for which the test set negatives are generated using Random Under-Sampling (RUS) with the same sampling ratio as used in training, thereby emulating a randomly-sampled negative set and explicitly probing potential dependence on a particular negative-selection scheme (see revised Figures 8–9). Note that for these revised analyses, we use multiple RUS ratios (1:4, 1:10, 1:20, 1:30) and two loss-weighting configurations (balanced versus dynamic weights) for both the Past and Dynamic Domains, reporting precision, recall, F1, ROC–AUC, and PR–AUC across all forecast horizons (6h–48h). The results show that TCG-Net remains stable and consistently competitive under these randomly-sampled negatives: For the Past Domain, performance trends are preserved across lead times, and the best configurations (notably RUS 1:4 with

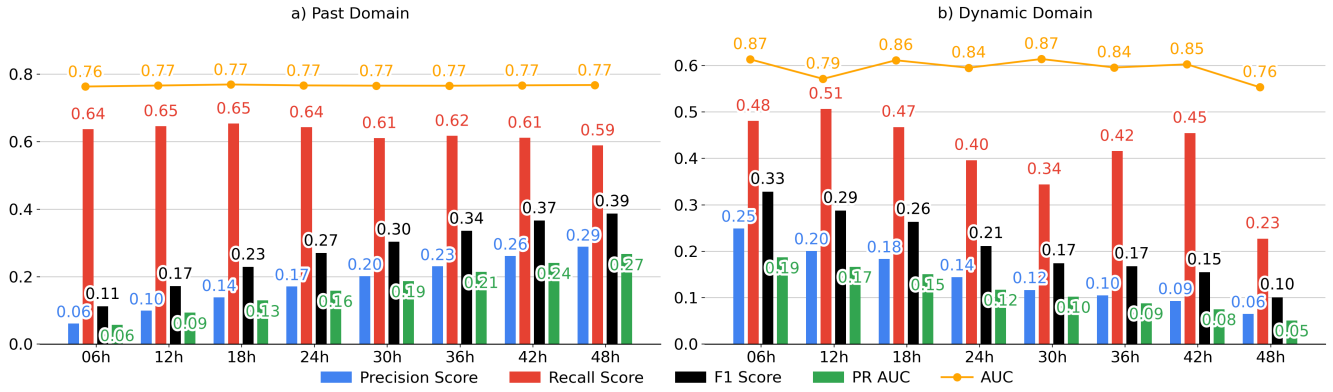


Figure 3: Overall performance of TCG-Net in terms of Precision (blue), Recall (red), F1 score (black), precision-recall area under the curve (PR-AUC, green) and area-under-the-curve ROC (AUC-ROC, yellow) for the TCG prediction on a) Past Domain, and b) Dynamic Domain.

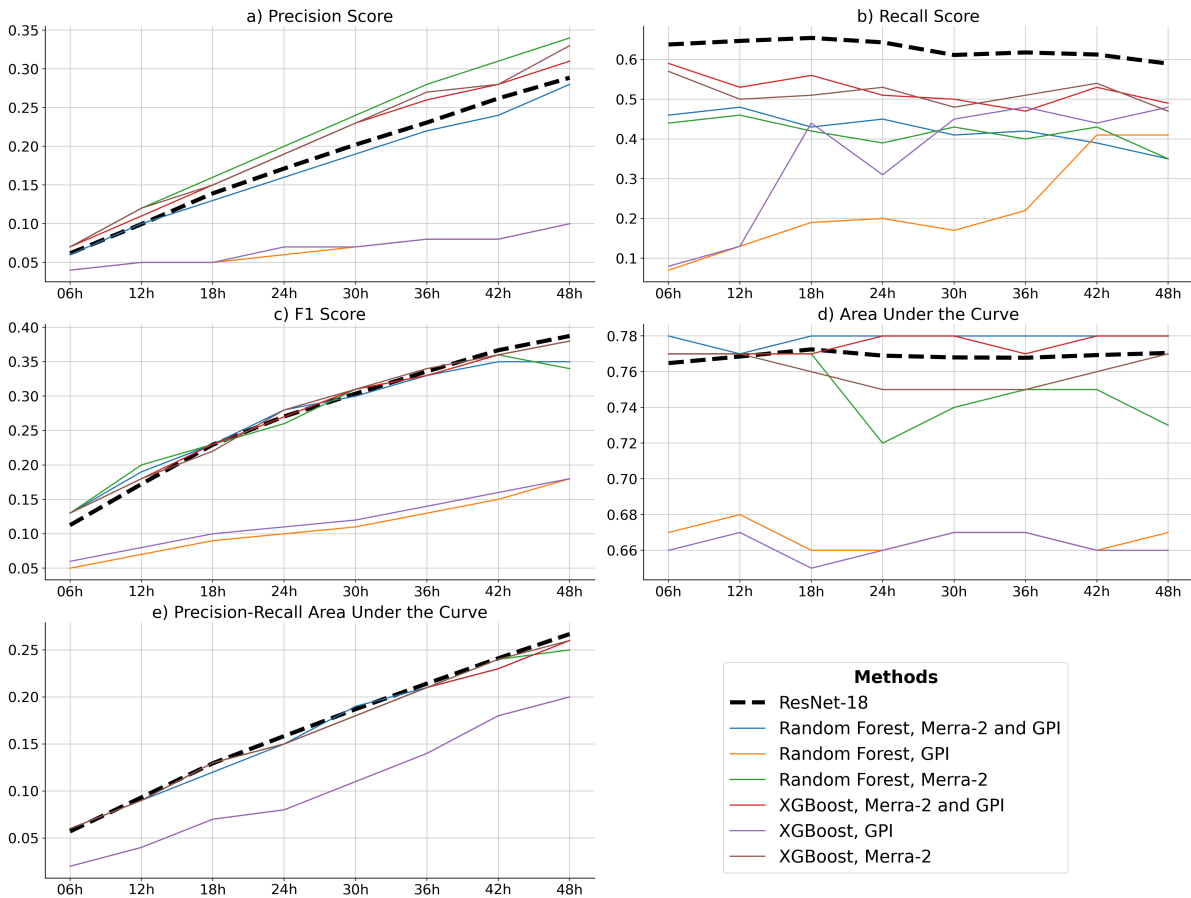


Figure 4: Comparison of the model performance for the **Past Domain** task between TCG-Net and traditional classification models, i.e., XGboost and random forest, which are enhanced by including the climatological Genesis Potential Index (GPI). The thick dashed line denotes the performance of TCG-Net with ResNet-18 backbone.

dynamic weighting) maintain strong precision/F1 together with high ROC-AUC (0.78) and PR-AUC (0.85). For the Dynamic Domain method for which the task is intrinsically harder and non-stationarity is stronger, the model performance degrades with increasing lead time for all configurations, yet the relative ranking of methods and the overall behavior remain consistent. This indicates that TCG-Net's skill is not an artifact of a specific negative set construction. These updated results are now included in the revised manuscript, strengthening the evidence that our proposed DL approach could generalize under alternative and randomized negative-sampling conditions.

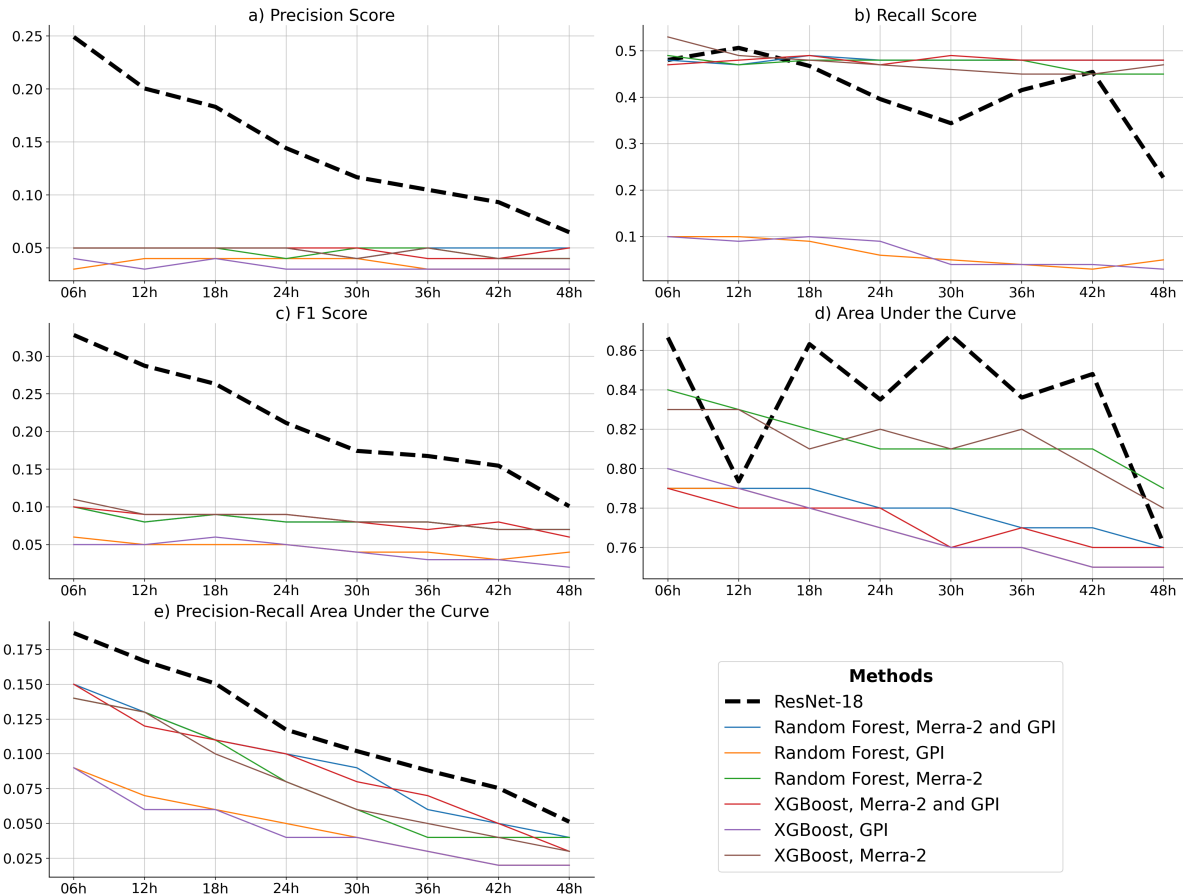


Figure 5: Comparison of the model performance for the **Dynamic Domain** task between TCG-Net and traditional classification models, i.e., XGboost and random forest, which are enhanced by including the climatological Genesis Potential Index (GPI). The thick dashed line denotes the performance of TCG-Net with ResNet-18 backbone.

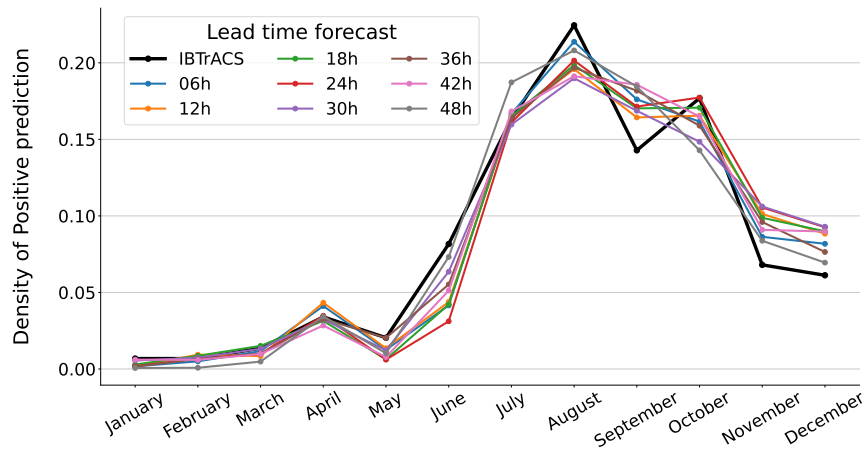


Figure 6: Monthly distribution of TCG frequency detected in the WNP basin from the test data (2017-2022), using the best-tuned ResNet-18 model for the DD strategy with data enrichment windows from 6 to 48 hours. The black solid curve denotes the TCG frequency obtained from the best track during the same time period.

## 2. Comparison with Traditional Method

*This study lacks a comparison with traditional approaches, so it remains unclear what the benefits and trade-offs of training TCG-Net actually are. For example, how would a Random Forest/LightGBM model perform? Alternatively, the authors should at least compare against a traditional index-based method, such as the classic*

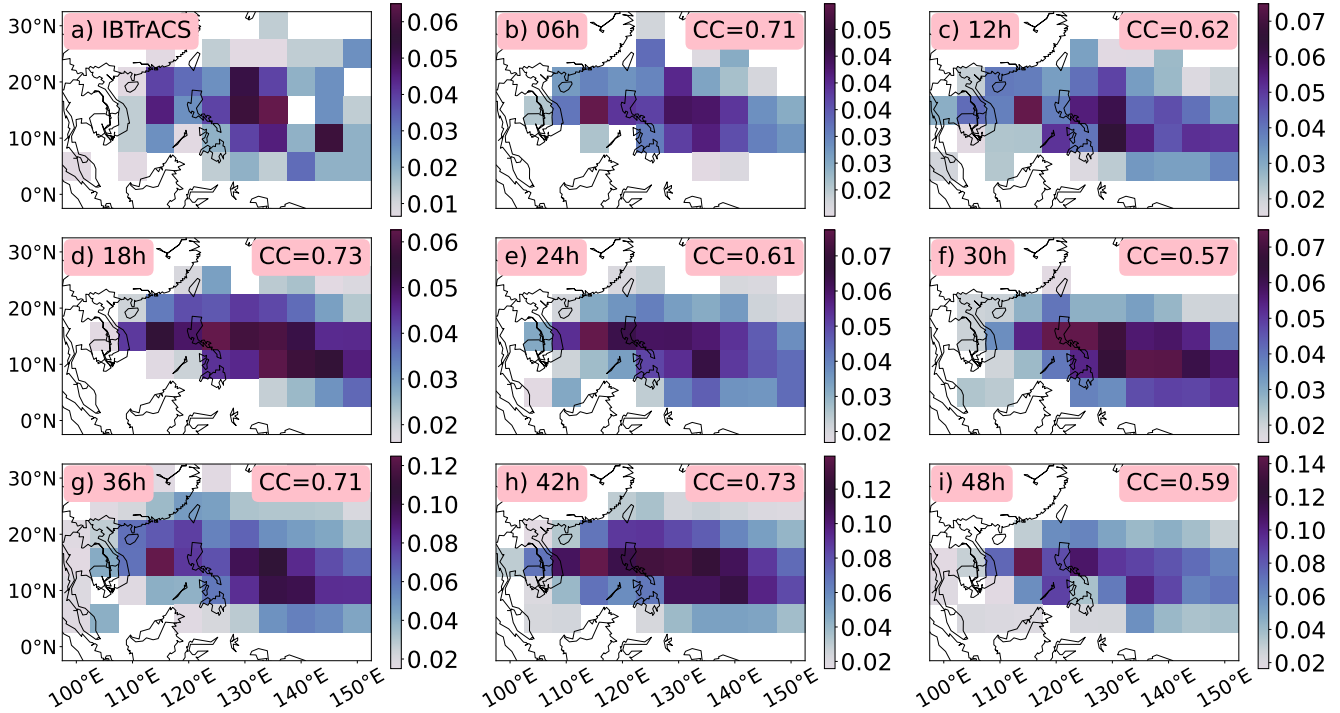


Figure 7: a) The spatial distribution of the observed TCG density (shaded) during the 2017-2022 period as obtained from the best track; (b)-(i) 5-year average of TCG probability prediction that is obtained from the ResNet-18 model with different data enrichment windows from 6-48 hours during the same test period as in (a). Note that different shading scales are used for different data enrichment windows so that one can better see the contrast between the areas of maximum probability for TCG predicted by the ResNet-18 model.

*Genesis Potential Index (GPI), which is not even suitable for 6-hour interval prediction.*

**Authors' response:** Thank you very much. Given your suggestion regarding benchmarking against traditional approaches, we have conducted in this revision additional experiments comparing TCG-Net with several widely used machine learning baselines including Random Forest and XGBoost, along with different combinations of MERRA-2 only, GPI only, and GPI as an additional feature of TCG-Net (revised Figures 4-5). Note that we also explicitly included the traditional GPI-based predictors to assess the value of this physically-derived index. The new results show that while tree-based models achieve reasonable recall at short lead times, their precision, F1 score, and especially PR-AUC degrade substantially with increasing forecast lead time, particularly for the Dynamic Domain method. In contrast, TCG-Net with ResNet-18 consistently maintains higher F1 and PR-AUC across all lead times (6-48 hours), indicating stronger robustness under class imbalance and temporal variability. Moreover, GPI-only models exhibit limited predictive skill, especially for 6-hour interval forecasting, confirming that purely GPI-based approaches are not sufficient for TCG prediction. Overall, the extended benchmarking demonstrates that the performance gains of TCG-Net are systematic rather than incidental, and that its training complexity yields measurable improvements over both classical machine learning and index-based methods.

### 3. The uncertainty of TC track data

*In this study, the authors define the first reported position in the TC best-track dataset as the genesis location. It would be important to conduct a sensitivity test to assess how the model's performance changes if the genesis time/location is shifted slightly (e.g., by  $\pm 1 \sim 3$  hours or  $\pm 1$  degree). Since this study applies hourly reanalysis dataset rather than monthly dataset to construct DL models, such analysis is necessary.*

**Authors' response:** Your comment directly addresses a central challenge in TCG climatology reconstruction, for which we fully agree with. In fact, the timing and location of a TCG event derived from best-track datasets can vary substantially in both time and space in some cases, especially during the early INVEST stage when it is often

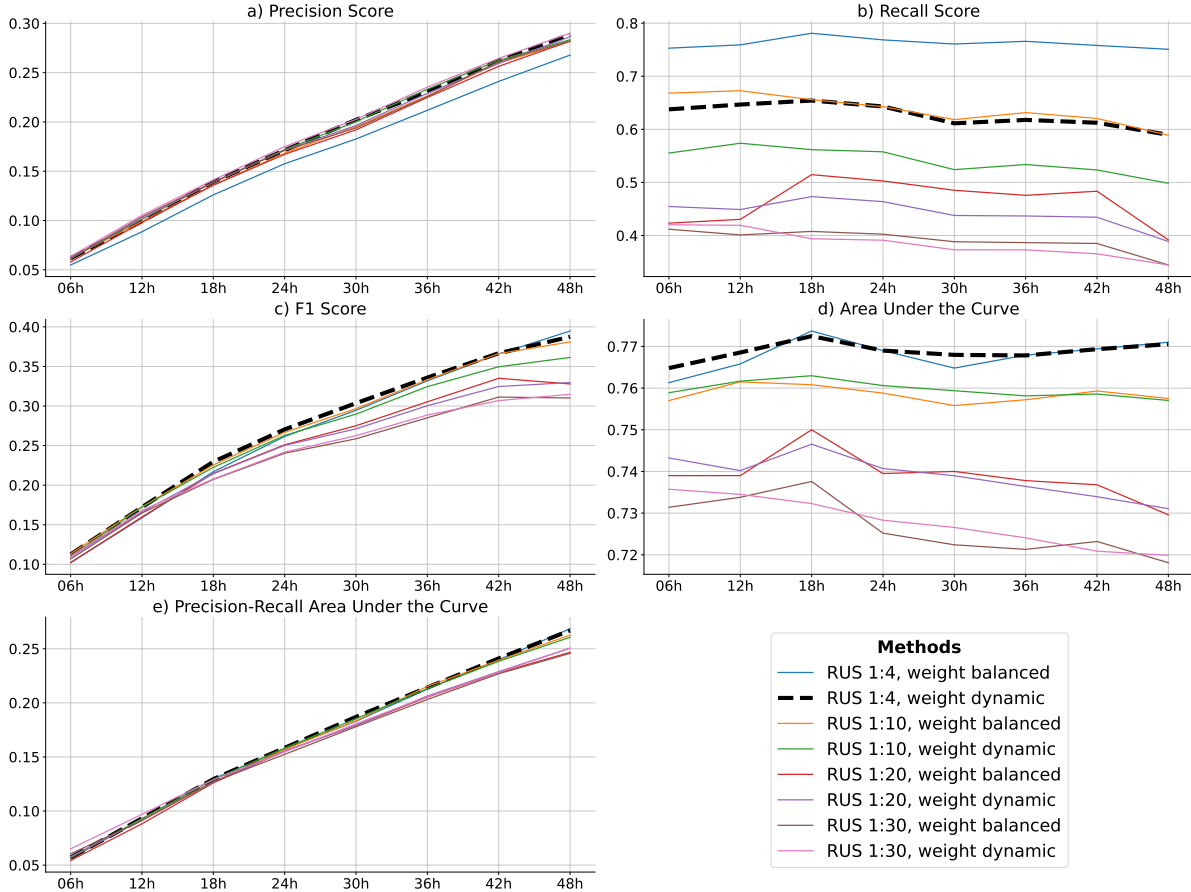


Figure 8: (a) The precision score  $P$  for the TCG prediction of TCG-Net using a range of the RUS ratio and class weight (solid colors) for the **Past Domain** task and the same sampling ratio for both the training and test sets; (b)-(c) similar to (a) but for the  $R$  and F1 scores, respectively. The dashed black line denotes the reference obtained from our best-tuned model. Note that *weight balanced* assigns fixed importance to each class based on frequency whilst *weight dynamics* adaptively adjusts sample or class importance.

difficult to identify the exact onset time of TCG. This uncertainty is exactly why we introduce a data-enrichment strategy that extends beyond the single time stamp obtained from the best-track record. Specifically, as described in the Methods and Results sections, we label positive TCG events using a retrospective temporal window of up to 48 hours prior to the best-track TCG timing in all of our analyses. In addition, spatial uncertainty in TCG location is explicitly considered by defining a relatively large positive-label domain ( $\sim 18 \times 18$  degrees), thus making sure that the favorable environment for a TCG event is fully captured. These inherent temporal and spatial uncertainties are also the reasons why we use sliding-window approaches when reconstructing TCG climatology, as illustrated in Figures 6–7.

In response to your comment, we have expanded the discussion of these issues in the revised manuscript and hope this more clearly conveys the rationale behind our methodological choices and the interpretation of the resulting climatological patterns.

### 5. XAI method

In Section 4.3, the authors conduct sensitivity analyses by comparing different sets of input reanalysis variables. However, for modern deep-learning models it is straightforward to include gradient-based XAI diagnostics. For example, the authors could compute saliency maps (or related methods such as Integrated Gradients/Grad-CAM) to visualize which regions and which variables most influence the genesis probability, and then quantify the relative importance of each variable based on these attributions. This would provide a more direct and model-consistent assessment of variable importance than input-subset sensitivity tests alone.

**Authors’ response:** Agree. Gradient-based XAI methods are indeed very useful for analyzing the role of different variables in probability maps for specific case studies. In our previous work (Nguyen and Kieu, 2024), for example,

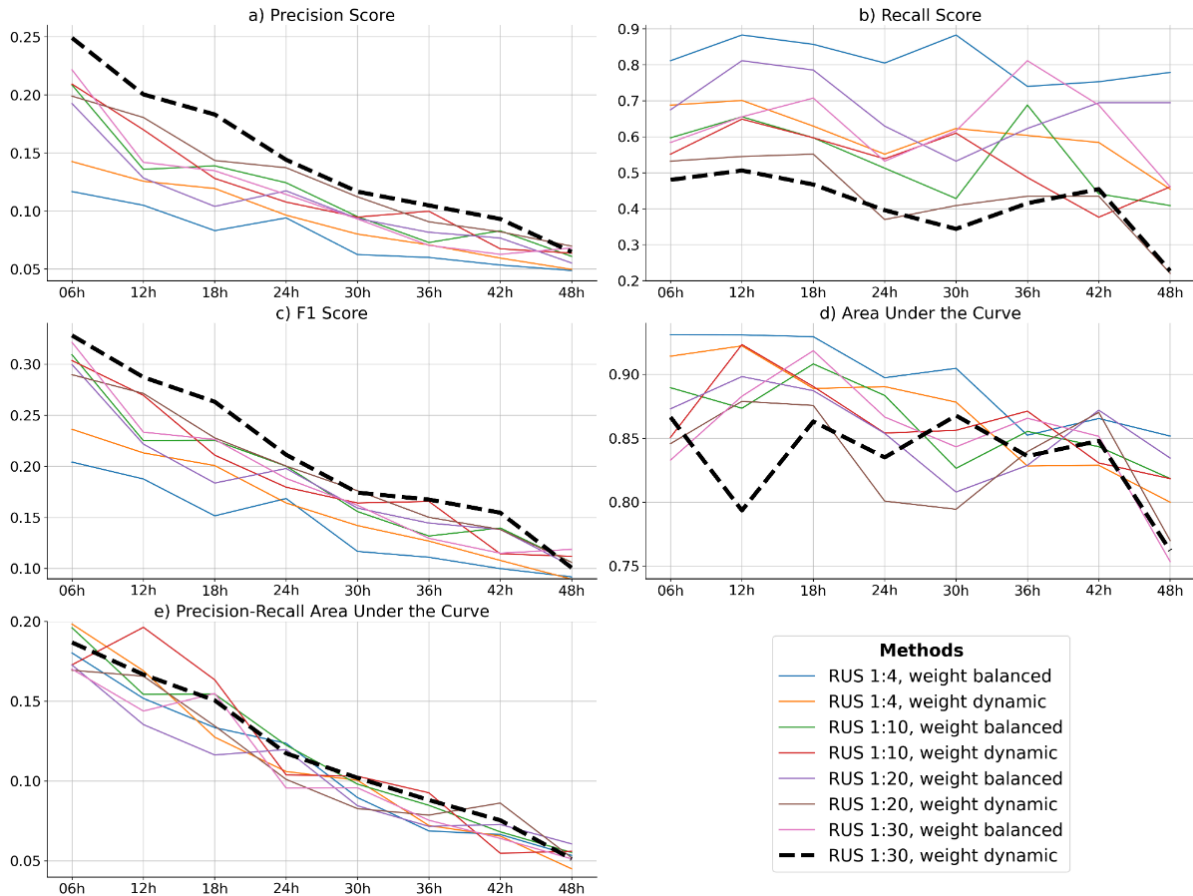


Figure 9: (a) The precision score  $P$  for the TCG prediction of TCG-Net using a range of the RUS ratio and class weight (solid colors) for the **Dynamic Domain** task and the same sampling ratio for both the training and test sets; (b)-(c) similar to (a) but for the  $R$  and  $F1$  scores, respectively. The dashed black line denotes the reference obtained from our best-tuned model. Note that *weight balanced* assigns fixed importance to each class based on frequency whilst *weight dynamics* adaptively adjusts sample or class importance.

the integrated-gradients method was applied to identify where information contributing to a predicted TCG event was drawn from. In this study, our focus focuses however more on reconstructing TCG climatology rather than analyzing individual events. Applying a gradient-based XAI approach for each case during the reconstruction would therefore be overwhelming. Our main aim in this study is to design a DL model that can represent the climate mean robustly, rather than individual realizations as typically emphasized in weather-forecast applications. As such, we adopt in this study an XAI framework that emphasizes sensitivity analyses with respect to different input channels, model parameters, and seasonality shown in Figures 6, 12–15, and Table 4, which can be naturally interpreted from a climatological perspective. These sensitivity analyses provide insight into how variations in each input channel influence the reconstructed TCG climatology as a whole, and so they are a form of XAI by nature.

In response to your comment, we have now included additional statistics on the ranking of each input variable based on its weights, which is averaged from the first CNN block of the RestNet-18 model (Figures 10–12). This ranking offers further information on the relative contribution of each channel to the overall TCG climatology reconstruction as expected.

*Minor Comments*

*P4 L95: The rationale for selecting MERRA-2 needs to be strengthened. First, TC climatology reconstruction does not necessarily require very high resolution, since the ML models in this study are not tracking individual TCs. Do the authors analyze how fast TC moves when it is generated, and how uncertain the first location of TC best track represents TCG?*

*Second, in most ML (especially DL) workflows, the original data are typically converted into standardized analysis-ready formats (e.g., zarr), so data format alone is not a compelling reason to prefer a particular reanal-*

Table 4: List of features selected using the feature engineering and feature ranking filter approach as obtained for each labeling strategy during the training period.

| ID | Name of Features | Past Domain         |   | Dynamic Domain      |   |
|----|------------------|---------------------|---|---------------------|---|
|    |                  | Feature Engineering | Feature Ranking                         | Feature Engineering | Feature Ranking   |
| 1  | QL               |                     | 400, 700, 825, 900, 950                 |                     | 100, 1000, 150, 200, 300, 400, 500, 600, 700, 800, 875, 900, 950, 975 |
| 2  | H                | 500                 | 200, 925                                | 500                 | 100, 550, 950   |
| 3  | QI               |                     | 250, 450, 600, 800, 900, 925, 950, 1000 |                     | 100, 1000, 150, 500, 600, 700, 900                                    |
| 4  | OMEGA            | 500                 | 450, 875                                | 500                 | 100, 150, 250, 600, 925, 1000   |
| 5  | T                | 500, 900            | 725                                     | 500, 900            | 150, 200, 900   |
| 6  | U                | 200, 800            | 825, 1000                               | 200, 800            | 1000, 550, 200  |
| 7  | V                | 200, 800            | 150, 550                                | 200, 800            | 1000, 600, 400, 150, 100  |
| 8  | RH               | 750                 | 950                                     | 750                 | 100, 200, 400, 700, 825, 875, 925, 1000                               |
| 9  | QV               |                     |   |                     | 100, 150, 900   |
| 10 | VOR              | 200, 700, 900       |   | 200, 700, 900       |   |
| 11 | DIV              | 200                 |   | 200                 |   |

*ysis. Please provide additional, science-based justification—for example, whether MERRA-2 has demonstrated advantages over other reanalyses in representing key genesis-relevant environmental fields such as vertical wind shear and low- to mid-level humidity.*

**Authors’ response:** Thank you. We wish to take this opportunity to clarify that our use of MERRA-2 at 0.5-degree spatial resolution is primarily motivated by its similar resolution to most of the global climate projection products from current CMIP5/CMIP6. Training the deep learning (DL) model at this resolution, thus, facilitates subsequent fine-tuning with other climate datasets, which is our further aim after this study. In fact, our ongoing extension of this work that aims at reconstructing TCG during the pre-satellite era follows this exact strategy. That is, our DL model is first trained using MERRA-2 data and then further fine-tuned with ERA5. A similar approach could be adopted using other reanalysis datasets, such as JMA, CFS, or NCEP-FNL, all of which are available at 0.5-degree resolution. To the best of our knowledge, there is currently no comprehensive study comparing TCG climatology across these reanalysis products. As such, any of these datasets is equally valuable for DL model development. This consideration underlines our decision to present a DL model development using the MERRA-2 dataset in this study. These discussions have been included in this revision per your comment.

Regarding the Zarr format, we could see that its use is indeed becoming increasingly popular within the atmospheric community. At present, most existing climate datasets are distributed in NetCDF or GRIB1/2 formats. For this reason, we chose to output our preprocessed data in the same original NetCDF format to maintain a consistent and unified framework across all I/O stages (and to minimize use of external libraries). Incorporating additional data formats would require only minimal modifications to our workflow, as users can readily extend the preprocessing component to support alternative formats and generate a common intermediate data representation.

For the final note on TC tracking, we would like to mention that our main aim of this study is for TCG, which concerns only the very first location when a TC forms, instead of the entire TC track. Thus, we do not have any analyses for TC movement, which will require a new component for TC vortex tracking beyond the scope of this work.

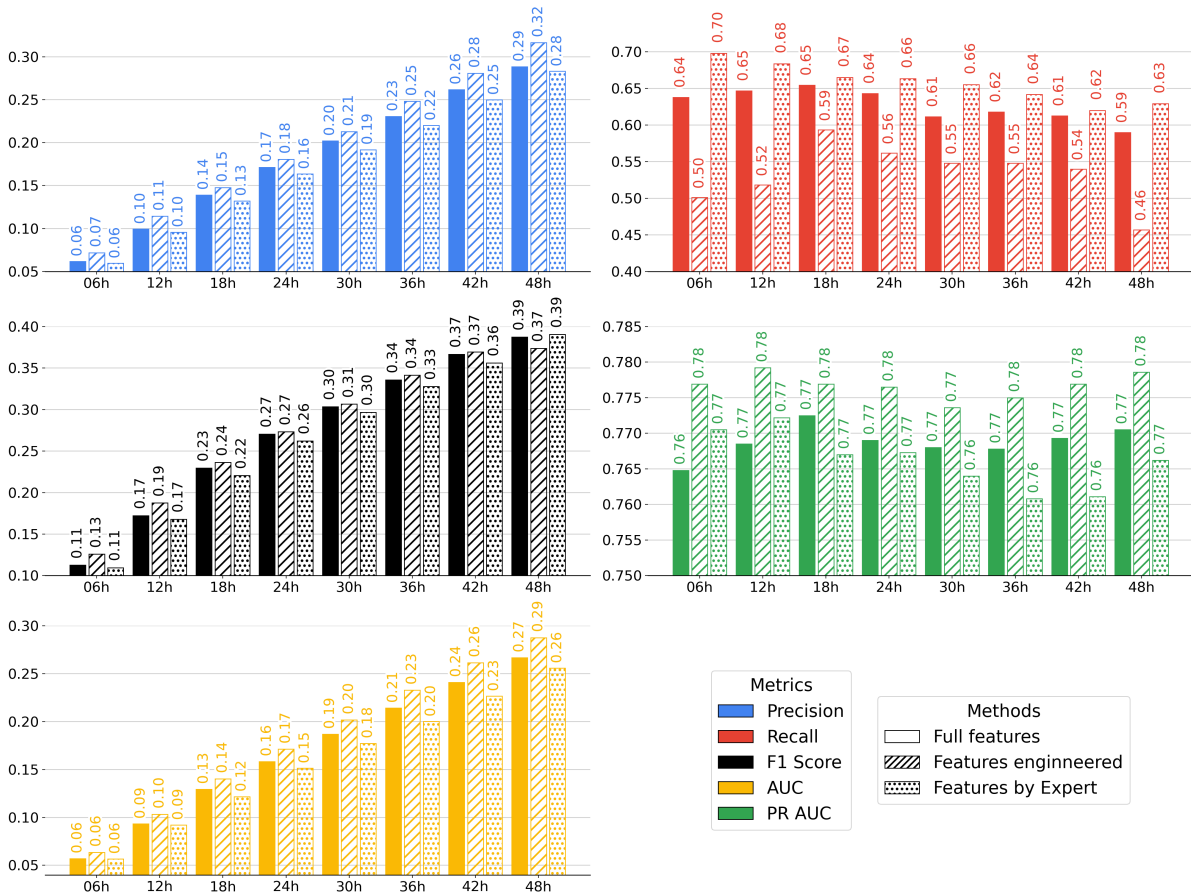


Figure 12: Evaluation of the model performance for several different feature selection methods including all features (solid color columns), 13 selected features based on feature engineering in (Nguyen and Kieu, 2024) (striped columns), and feature ranking of top 10% (dotted columns) using the **Past Domain** task.

## Response to Reviewer 2

*This manuscript presents TCG-Net, a deep-learning-based framework for the reconstruction of the tropical cyclone genesis (TCG) distribution over the western North Pacific (WNP). The authors introduce two task-specific labelling strategies, combined with temporal feature enrichment and imbalance-aware training, to extract both seasonal and spatial characteristics of TCG directly from MERRA-2 reanalysis data. The effort put forth by the authors is commendable, but the manuscript requires significant revision. In particular, clearer clarification is needed on the applicability of TCG-Net, the definition and climatological representation of TCG, the claimed physical novelty. Before considering this paper for publication, I have several concerns and suggestions outlined below.*

**Authors' response:** We are grateful to Reviewer 2 for your thorough evaluation and detailed suggestions. In this revision, we have carefully considered all of your comments and have made substantial changes to address the concerns you raised, as detailed in our point-by-point responses below. All revisions are indicated in the tracked-changes version for your convenience. We hope that the revised manuscript meets your expectations and will be acceptable to you.

*1. TCG-Net is an application-oriented framework, yet its practical strengths and limitations relative to existing approaches are not sufficiently clarified. The authors claim that vortex tracking methods face challenges with coarse-resolution climate models ( $> 0.5^\circ$ ) and TCG-Net therefore serve as a complementary tool. In fact, several traditional detection algorithms has been developed for coarse-resolution datasets, such as OWZP and TRACK*

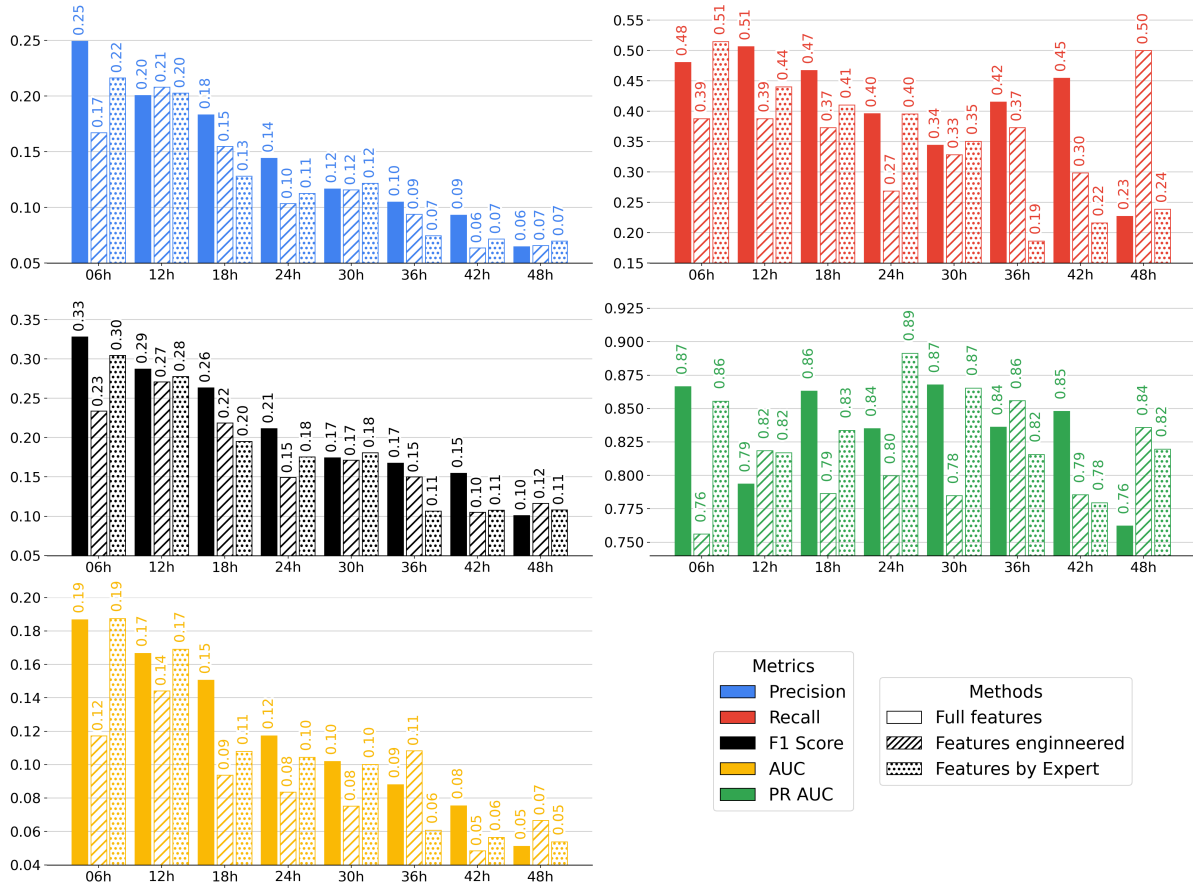


Figure 13: Evaluation of the model performance for several different feature selection methods including all features (solid color columns), 13 selected features based on feature engineering in (Nguyen and Kieu, 2024) (striped columns), and feature ranking of top 10% (dotted columns) using the **Dynamic Domain** task.

(Tory et al., 2013; Hodges et al., 2017), which exhibit reliable performance in reproducing the climatological distribution of TCG in both reanalysis datasets and climate models (Bell et al., 2019; Bourdin et al., 2022). A comparison between TCG-Net and traditional detection algorithms would be valuable.

**Authors' response:** Thank you for this comment. After conducting this type of comparison, we feel that such a comparison would be challenging to implement. The primary difficulty arises from differences in spatial resolution, aim, and study designs. Recall that our DL reconstruction is applied to a  $0.5^\circ$  resolution, whereas previous studies that you suggested (e.g., Tory et al.) derived TC climatology at coarser resolutions ( $1^\circ \times 1^\circ$  or  $1.5^\circ \times 1.5^\circ$ ). Applying their methods directly to our dataset would require a complete re-tuning of the vortex-tracking parameters in the OWZP framework to accommodate the higher resolution. Such re-tuning would introduce an additional layer of uncertainty and make a direct comparison between their established algorithm and our DL-based approach subjective and potentially less meaningful. Similarly, the study by Hodges et al. focused on detecting TCs from reanalysis datasets rather than constructing TCG climatology. As a result, their vortex-tracking methodology is not directly comparable to the TCG climatology framework used in our study.

Beyond the few you suggested, we have been exploring other alternative vortex-tracking methods, such as TempestExtremes (Ullrich and Zarzycki 2017), but have not yet obtained stable results due to some technical issues. Figure 16 in this response shows an example of the TCG distribution when we apply vortex detection directly on MERRA-2 data that we obtain from TempestExtreme for a test period from 2017-2022. The result turns out to be very sensitive to several parameters such as TC lifetime, vortex strength, or the surface wind threshold, which all affect the first moment that a TC is recorded for the TCG reconstruction.

At present, we still do not have a systematic way to adjust these parameters for every test period or data enrichment

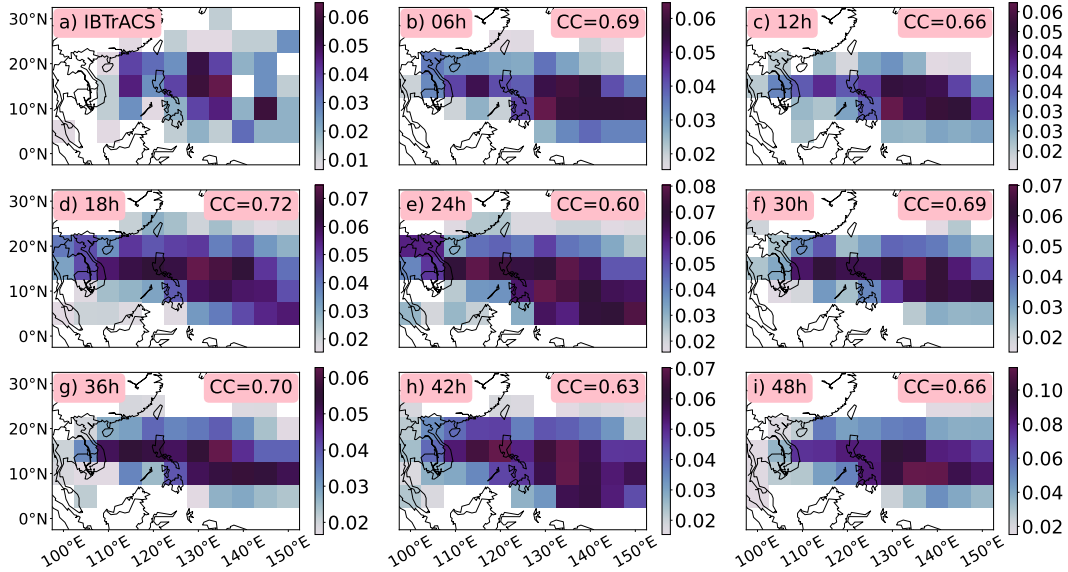


Figure 14: Similar to Fig. 7 but for the feature engineering approach.

window for this analysis. For these reasons, we have decided to restrict our analysis to a direct comparison between the DL-reconstructed TCG climatology and that derived from best-track data to minimize introducing more subjective analyses. We hope Reviewer 2 will agree with us on this approach.

*2. TCG-Net was not developed based on objective TC structural characteristics, but rather on large-scale environmental factors in present-day climates, similar to GPIs. As a result, the applicability of TCG-Net to future climate projections is uncertain, as the relative importance of critical environmental factors may change under warming conditions (Murakami and Wang, 2022).*

**Authors’ response:** This point is indeed central to our study. We would like to emphasize that the traditional approach to TCG climatology based on indices such as the Genesis Potential Index (GPI) and other well-established variables such as SST, CAPE, and vertical wind shear relies largely on empirical relationships or approximate theoretical frameworks. These commonly-cited environmental conditions for TCG have been derived from previous observations, field campaigns, and modeling studies. While they are statistically robust and physically meaningful, they are not universally sufficient to guarantee TCG across all basins. Each basin has distinct regional characteristics that may enhance, modulate, or suppress particular processes.

In this study, we introduce a DL framework that systematically evaluates a broad range of environmental variables and identifies the dominant controlling factors within a unified modeling architecture. As presented in Section 3, our results confirm several key variables that are consistent with existing physical understanding, while also revealing additional factors that appear to be specific to the WNP basin. These newly identified factors provide further insight into TCG behavior in the WNP and may offer added value for assessing potential future changes in TCG under a changing climate. Importantly, our approach is systematic and readily extendable to other basins, datasets, or climate regimes to identify the most relevant controlling mechanisms in each context.

In response to your comment, we have incorporated new analyses in which GPI is (1) used as a proxy for TCG and (2) included as an additional input channel to the DL model. The results (Figs. 4–5) show that GPI alone does not perform as effectively as DL reconstruction. On the other hand, including GPI as an independent input channel contributes little to overall model performance, likely because the environmental variables used to compute GPI are already included among our existing input channels. In this revision, we have updated Section 3 to provide additional discussion regarding the potential role and contribution of GPI, as you suggested.

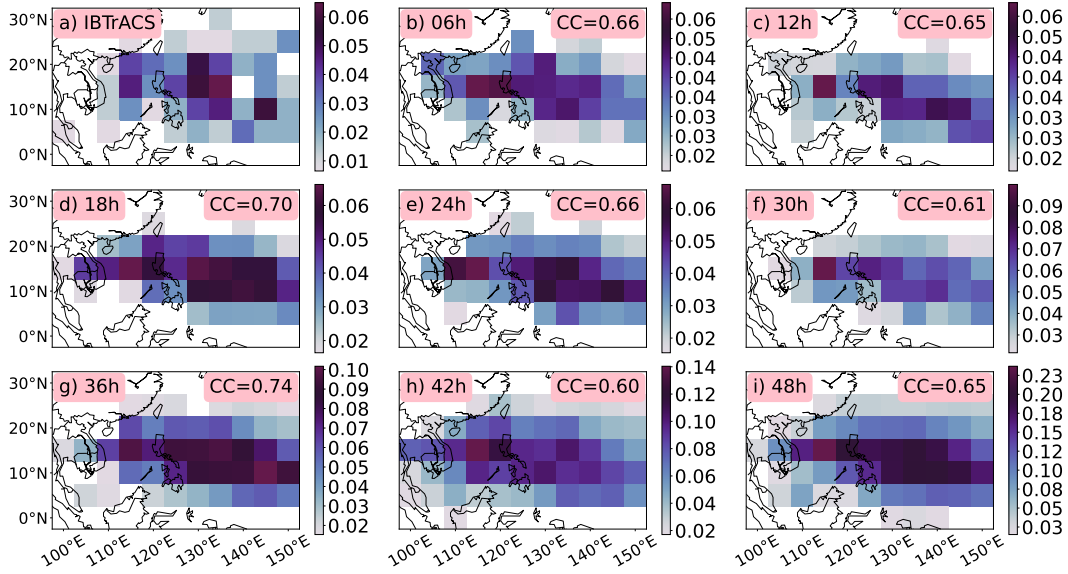


Figure 15: Similar to Fig. 7 but for the automatic feature ranking approach.

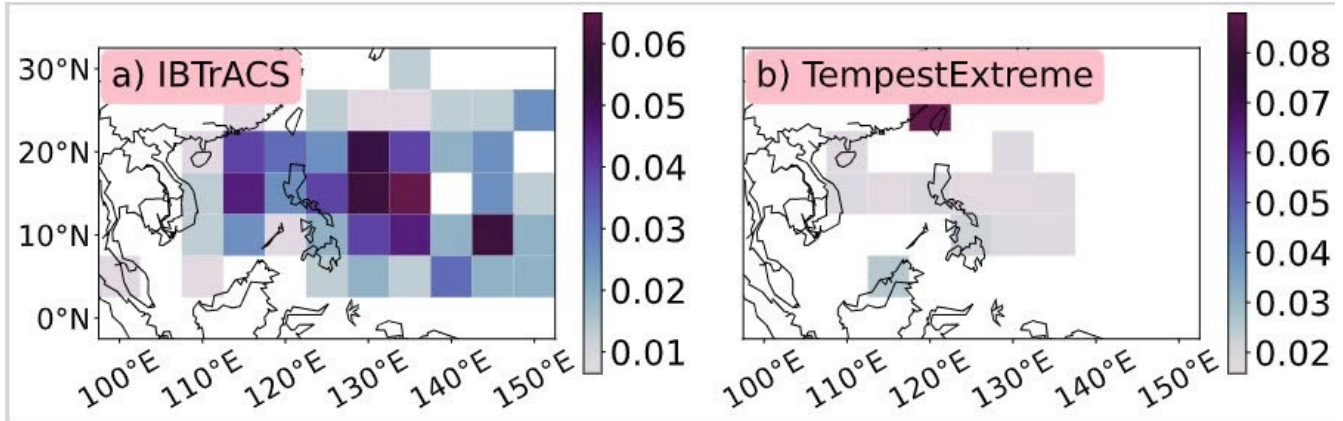


Figure 16: Spatial distribution of TCG density as obtained from a) best-track, and b) tempestExtreme vortex tracker for a test period 2017–2022, using the set of vortex tracking thresholds for the NCEP 0.5-degree resolution.

3. The authors use the TCG distribution derived from 2017–2022 to represent “TCG climatology,” which may not be appropriate. Given the strong interannual variability of TCG, a five-year period is generally insufficient to characterize climatological conditions.

**Authors’ response:** Your point is well taken. We acknowledge that this is a key limitation of our study, as a 5-year period cannot be considered fully representative of TCG climatology. The primary constraint preventing the use of a longer evaluation period is the availability of the MERRA-2 dataset (or any reanalysis datasets in general), which are most reliable from 1980 onward. With roughly 40 years of data in total, allocating 5 years for testing already represents a substantial fraction of the data record, leaving about 35 years for training. Extending the testing period further would reduce the amount of training data, potentially limiting the model’s ability to generalize across different climate modes. Our strategy therefore prioritizes maximizing the training dataset so that the model can robustly learn the environmental conditions conducive to TCG. While alternative approaches such as multiple 5-year leave-one-out experiments are possible, they still do not fully resolve this issue because each test climatology would still be constrained to a limited 5-year window.

We have added a discussion of this limitation in the revised manuscript. This challenge is indeed inherent to any DL-based approaches applied to relatively short climate records, and we believe it is important to clearly inform readers of this potential limitation when interpreting the results as you commented.

4. *In this study, TCG was defined as the first time that a TC was recorded in observations, which can be reasonable in a weather prediction context because it allows earlier detection of cyclogenesis. However, in climatological studies, TCG is more commonly defined when the storm intensity first reaches 35 kt in order to exclude weak or short-lived vortices (Klotzbach et al., 2022; Lai and Toumi, 2023). How sensitive the performance of TCG-Net may be to the definition of TCG?*

**Authors' response:** Yes, you are correct that the timing of a TCG event derived from best-track datasets can vary substantially, especially during the early INVEST stage or weak systems for which it is often difficult to identify the exact onset time of TCG. This uncertainty is exactly why we introduce a data-enrichment strategy that extends beyond the single time stamp provided by the best-track record in all of our analyses (i.e., the time window from 6-48 hrs in Figs 2-15). Specifically, as described in the Methods and Results sections, we label positive TCG events using a retrospective temporal window of up to 48 hours prior to the best-track record. This way, the inherent temporal uncertainties can be captured by our use of multiple data-enrichment windows when reconstructing TCG climatology. As illustrated in, e.g., Figures 6–7, the performance of our TCG-Net is quite robust for the enrichment up to 36 hr before the genesis timing recorded in the best track. Taking further data beyond 36-hr will reduce the performance of our TCG-net model. In response to your comment, we have expanded the discussion of these issues in the revised manuscript and hope this could better address the uncertainty issue that you raised.

5. *While the authors emphasizes novelty in terms of large-scale environmental drivers of TCG, the selected large-scale factors based on feature ranking are largely consistent with previous studies (Emanuel, 2010; Wang and Murakami, 2020). In this regard, the results appear to largely reproduce established findings rather than provide genuinely new physical insights, and the claimed level of innovation in this aspect may be overstated.*

**Authors' response:** This comment is somewhat related to your earlier comment (#2) that we responded above. One of the major points of our DL approach is its ability to identify and evaluate the role of various factors, which are further highlighted in the new analyses of factor ranking in this revision (Figs 12-13). As shown in these results, we not only confirm several key factors for TCG as identified in previous studies, but also uncovering several new factors relevant to our TCG reconstruction such as the cloud water content (QI/QL) or lower troposphere temperature, depending on the sampling strategies (see also Table 4). The fact that our DL could recover the previous factors gives us confidence that these new factors are indeed significant for TCG in the WNP basin. Note that these additional factors are basin dependent, and so there may emerge other factors in different basins. Regardless of the basin, our DL approach can help uncover within a consistent framework that we wish to present in this study. In this revision, we have further revised the feature importance section to better highlight the significance of our approach. We hope that this addition will offer readers deeper insights into the TCG process in the WNP basin, extending beyond the climatology-based reconstruction in previous studies.

6. *I was a little confused about domain chosen in this study. While a positive TCG label was defined as the square box of size  $18^\circ \times 18^\circ$  centered on the first recorded TCG location, the ResNet-18 model was applied on each  $5 \times 5$  box. It is therefore unclear how the labeling domain and the prediction domain were reconciled during training and evaluation.*

**Authors' response:** Thank you for your comment. Our previous discussion was unclear here. To clarify, the ResNet domain is  $18^\circ \times 18^\circ$  degrees, but we shift the domain center by 5 degrees each time when sliding the domain so we can cover the entire WNP basin. The TCG probability obtained from each domain is then assigned to the center of that domain, allowing us to construct a map with a resolution of 5 degrees. We have revised this description to ensure that our approach is clearly presented, as per your suggestion.

*Minor comments*

1. *L95: Another advantage of choosing MERRA-2 is that TC-related information has been assimilated into MERRA-2, whereas ERA5 does not include this (Gelaro et al., 2017).*

2. L114: *IBTrACS database compiles global TC tracks information from multiple agencies. Which agency's observations were utilized in this study?*
3. L179: *Change "DM" to "DD"*
4. L362-363: *The discrepancy between the DL model results and observations may not solely reflect limitations in DL model optimization, but could also arise from deficiencies in the ability of large-scale environmental factors to reproduce the seasonal variability of TCG (Menkes et al., 2012; Tippett et al., 2011).*
5. L375-381: *The pattern correlation coefficients with the observations are encouraged to quantify the performance of the DL model.*
6. L432: *The pressure levels described here is inconsistent with the Table 4*

**Authors' response:**

1. Thank you for pointing this out the benefit of MERRA-2 that we were not aware of. The work by Gelaro et al. (2017) has now been included in this revision.
2. We use the JTWC agency in the IBTrACS data, which is now mentioned explicitly in this revision.
3. This paragraph has been revised.
4. Totally agree. This session has been revised accordingly.
5. Pattern correlations for all maps have been now included in this revision.
6. This inconsistency has now been corrected.