

## Response to Reviewer 1

Overall, I find the study interesting and the results valuable. In particular, the successful implementation of trajectory calculations within a climate model framework is a notable strength of the paper, as such analyses are often challenging due to the limited temporal and spatial resolution of typical climate model output.

That said, I believe the presentation of the material could be substantially improved. My comments below are therefore mainly concerned with the clarity and structure of the presentation, rather than with the underlying scientific approach or results.

We would like to thank the reviewer for their helpful comments. Our responses are shown in blue, while the reviewer's questions are in black. We also appreciate the suggestion to improve the clarity and structure of the presentation. In our revised manuscript, we have carefully reviewed it again and incorporated the reviewers' comments to enhance readability and make the narrative as clear as possible for the reader.

### Major Comments

1. The descriptions in the manuscript are often quite lengthy and at times somewhat repetitive. For example, results are first presented for DJF warming events in the present and future, then DJF cooling events, followed by JJA warming and cooling events, and subsequently for the tropics. To enhance reader engagement, it might be helpful to introduce the figures once and then focus more on highlighting the key differences, while keeping the descriptions concise and emphasising the essential points.

**Response:** We have not changed the entire structure of the manuscript. We think the details are valuable and can be interesting for readers interested in specific types of events. Furthermore, the common structure makes it easy for readers not interested in specific details to skip some parts of the manuscript, and the highlights are summarised in section 4.

2. Playing the devil's advocate: How physically meaningful are the computed trajectories within the atmospheric boundary layer and especially in the tropics, where turbulent mixing is intense? I think a brief discussion on the limitations or uncertainties associated with the trajectories in these regions would strengthen the study.

**Response:** Thank you for this insightful comment. We agree that discussing the limitations of trajectories in the (tropical) boundary layer has strengthened the study.

We have added a brief discussion acknowledging that:

Trajectory calculations rely on resolved wind fields and therefore cannot capture the full spectrum of convective and turbulent motions. This limitation is especially pronounced in the tropics, where intense moist convection introduces significant uncertainty into the calculated transport pathways (Bao & Stevens, 2021; Bergman & Sardeshmukh, 2004; Stohl, 1998).

## Minor Comments

L15/16: The mention of a “clear dipole pattern” is somewhat confusing to me. You mention a “clear dipole pattern”, but then for JJA, the pattern does not clearly take a dipole form. Consider rephrasing.

**Response:** We have rephrased that sentence as: "The projected changes in (extreme) DTDT variations display a seasonal contrast: weakening in mid- to high latitudes and intensification in the tropics during December–February (DJF), while during June–August (JJA), tropical intensification is more widespread, and only some extratropical locations experience reductions in DTDT variations”.

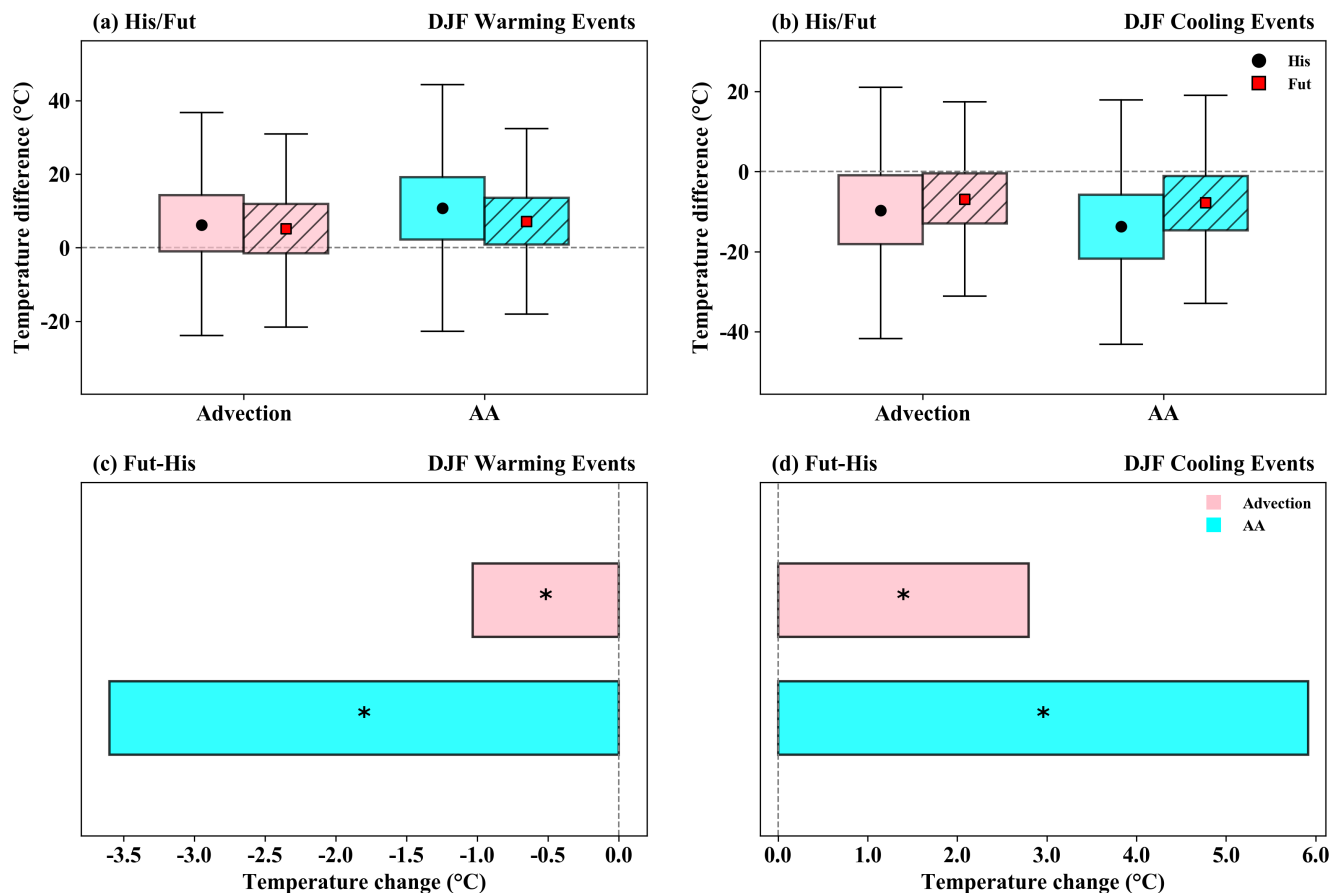
L18: “only” instead of “also”?

**Response:** We have changed this.

L19: I think you should be more careful here when writing “due to Arctic Amplification”.

**Response:** We have added some results based on this as per the suggestion of Reviewer 2.

Our trajectory analysis reveals a clear fingerprint of Arctic amplification: air originating from the Arctic experiences consistently larger future warming than air from lower latitudes (Figure 5, main paper). This reflects a weakening of the meridional temperature gradient. Furthermore, as shown in Figure R6 (response to Reviewer 2), air masses on both day  $t-1$  and day  $t$  are substantially modified by the warming climate. The ratio of local temperature change to global warming exceeds 1 on all days ( $-72$  h to 0 h), and is above 1.4 for both day  $t-1$  for warming and day  $t$  for cooling events, satisfying the quantitative definition of Arctic amplification. This confirms that air masses participating in these extreme DTDT events carry the signature of Arctic amplification, as they originate from regions warming at a rate significantly faster than the global average.



**Figure S1.** (a,b) Box plots comparing the contribution of advection and the estimated effect of Arctic Amplification (AA) for (a) warming events and (b) cooling events in North America, on a 3-day time scale. For each process, the left box (solid edge) represents the historical climate, and the right box (hatched) represents the future climate. The box spans the 25th to 75th percentiles; the black dot marks the historical mean; the red square marks the future mean; whiskers extend to 1.5 times the interquartile range. The effect of AA is estimated by comparing the seasonal-mean projected temperature change at the origins (at -72 h) of the historical backward trajectories for days  $t$  and  $t-1$ . (c,d) Stacked horizontal bars showing the mean change (future minus historical) for advection and AA for (c) warming and (d) cooling events. The \* indicates that the difference is statistically significant at the 95 % confidence level based on a t-test.

We isolated the Arctic Amplification (AA) contribution to extreme temperature events by comparing the seasonal-mean projected temperature change at the origins (at -72 h) of the historical backward trajectories between day  $t$  and day  $t-1$  in Figure S1. In this way, only the warming at the origins (which differs between  $t$  and  $t-1$  due to the typically different latitudes of air parcels at -72h) is considered under the assumption that the circulation is unchanged. The magnitudes of the AA effects are more pronounced than those from the changes in advection determined from equation 2. Specifically, for warming events, the change in advection (-1.1 °C) is smaller in magnitude than the change in AA (-3.5 °C), while for cooling events, the change in advection (+2.8 °C) is also smaller than the change in AA (+5.9 °C). The effect of AA explains the weakening of the advective contribution to both warming and cooling events, but the magnitude of the AA effect is mitigated by circulation changes. Also, the event-to-event variability (boxes and whiskers in Figure S1) is very similar between advection and estimated AA contributions.

We have added the Figure to the supplementary material and discussed this result in the main paper.

L37: “imperative” is a very strong word. I would be a bit more moderate here.

**Response:** We have changed this. Now the sentence reads “Therefore, studying DTDT changes and their extremes in a warming climate is important”.

L56: “for the past” instead of “in the past”?

**Response:** We have changed this.

L63: I would avoid citing Mayer (2025) when discussing the importance of diabatic heating, as Mayer (2025) emphasises the role of advection in temperature extremes rather than adiabatic or diabatic processes.

**Response:** We have removed this.

L69: “process understanding” instead of “processes understanding”?

**Response:** We have changed this.

L95: “use” instead of “utilise”

**Response:** We have changed this.

L106: In its current position, the formula does not appear to be well integrated into the text. The same applies to Eq. (2).

**Response:** In the revised manuscript, we have restructured the surrounding text to properly introduce and discuss each equation before it appears.

L110: Why are these seasons “key seasons”? Rather explain or omit the “key”.

**Response:** We have removed this and improved the sentence.

L124-127: I would omit the lengthy description of all the individual grid points in the supplement.

**Response:** We agree and have omitted the lengthy description of individual grid points and added it as a supplement.

Eq. (2): It might be helpful to write out the integrals explicitly to clarify exactly which terms are being computed.

**Response:** As the explicit integral forms are the same as those shown in Part I, we have added a cross-reference in the text to guide readers to the complete expressions.

We have added this: "Here, the DTDT change ( $\delta_T^0$ ) is decomposed into three contributing factors in Eq. (2), with full integral expressions calculated as in Part I (equations A4-6 in (Hamal & Pfahl, 2025))."

Eq. (2): Why do you accumulate over 3 days? Is there a physical reason? Have you tested the sensitivity of your results to other accumulation periods?

**Response:** We use LAGRANTO to calculate 10-day backward trajectories initialised at 18 UTC on both the preceding day ( $t-1$ ) and the event day ( $t$ ), at 10, 30, 50, and 100 hPa above the surface, for the corresponding grid boxes. The first part of our study, as well as previous studies, suggests that extremes typically develop on a timescale of 2–3 days (Bieli et al., 2015; Röthlisberger & Papritz, 2023). Therefore, we focus on 3-day backward trajectories for our analysis. For tropical locations, the time scale is further reduced to 1-day, since the relevant processes appear to occur on shorter time scales (see section 3.2.2 and Fig. 10). This has been discussed in Part I of our study.

L134-136: It is not clear what is meant by “mean temperature difference” or “mean adiabatic compression.” Consider clarifying the meaning of “mean” here.

**Response:** We have clarified that "mean" here refers to the trajectory-average value of each process over the 3-day back trajectory.

L147/148: There seems to be a contradiction: first, the results are said to fit ERA5 “in many regions,” then to deviate “in large parts.” Consider clarifying this.

**Response:** In the original manuscript, we used a different method to compare ERA5 and CESM. Following Reviewer 2's recommendation, we have now employed bootstrapping with the FDR test for this analysis. The text has been updated accordingly in the revised manuscript.

Figure 1, 2, 3, etc.: I think it would be helpful to use white color for small deviations around 0.

**Response:** The newly introduced significance test based on bootstrapping indicates that relatively small changes can also be statistically significant (see our responses to Reviewer 2 for more details). To prevent these signals from being masked, we did not adjust the color scale.

L162: The pattern does not appear as “distinct” to me.

**Response:** We have removed this.

L166/167/177: Phrases like “changes are driven by”, “increases due to”, or “influence” imply causality to me. Since the decomposition (into  $\sigma_T$  and autocorrelation) is “just” descriptive, I think you avoid implying causality.

**Response:** We thank the reviewer for this important distinction. We have revised the text to avoid causal language (e.g., "driven by," "due to") and instead use descriptive terms (e.g., "associated with," "contributed to by," or "coincides with") that align with the decomposition framework.

L179/180: The statements about Chile are contradictory: first mentioning it as an exception, then saying “(apart from Chile).” Consider clarifying.

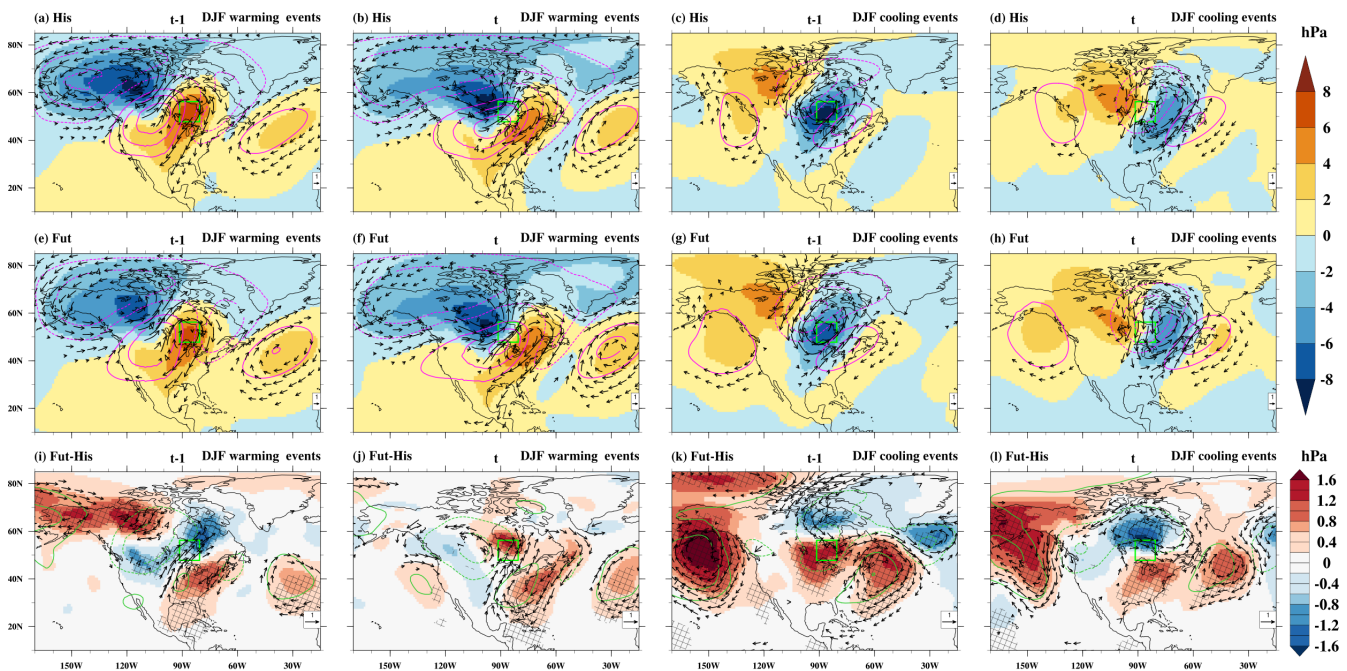
**Response:** We have made this clear based on the revised Figure.

Figure S3: The description is very detailed, but the figure is in the supplement. Consider either shortening the description or moving the figure to the main text.

**Response:** The description largely follows that of Figures 1 and 2 (in the main paper), with the only difference being the statistical analysis. We have therefore shortened the paragraph in the revised manuscript.

Figure 3/4/7/12/13: The figures are very small, which makes it difficult for the reader to fully appreciate their content.

**Response:** Thank you for the suggestion. We have now slightly enlarged the composite Figures, changed the stippling to more visible cross-hatching, and reduced the spacing between the geopotential height and wind panels, as shown in the revised Figure S2 below.



**Figure S2.** Composite of sea level pressure anomalies (hPa, color shading), wind anomalies at 850 hPa ( $\text{m s}^{-1}$ , vectors), and geopotential height anomalies at 500 hPa (gpm, magenta and darkgreen contours) relative to the seasonal mean on the (a, e, i, c, g, k) previous day (t-1) and (b, f, j, d, h, l) the event day (t) of the warming (a-b, e-f and i-j) and cooling (c-d, g-h and k-l) events during December-February (DJF) in (a-d) historical climate (His), (e-h) future climate (Fut), and (i-l) projected future changes (Fut-His) at a selected grid box in North America (green box). Note that, in (a-h), wind vector anomalies  $\geq 2 \text{ m s}^{-1}$  and in (i-l), wind vector difference anomalies  $\geq 0.5 \text{ m s}^{-1}$  are plotted. The dotted and bold contours indicate negative and positive geopotential height anomalies, respectively. Additionally, the cross-hatched area indicates where the ensemble mean of sea level pressure differences exceeds the 95% confidence threshold based on a t-test.

L208: omit the “future” as projected already implies future?

**Response:** We agree and have removed "future" since "projected" inherently refers to future conditions.

L239: Out of curiosity: Do you have an idea why the contribution of the diabatic heating is larger in the CESM-LE compared to ERA5? Could this relate to vertical resolution?

**Response:** Thank you for your curiosity. We have added the new description and analysis in the revised manuscript.

Approximately 40–70% of global land regions exhibit systematic biases in daily temperature and related metrics simulated by CESM-LE compared with ERA5 (Figures R2–R3 in our response to Reviewer 2). While the model captures large-scale spatial patterns (cf. Part I), there are notable discrepancies in the magnitude of anomalies (Figure R3) and in the contributions of underlying physical processes. These biases include an overestimation of DTD extremes, most pronounced in the mid-to-high latitudes (e.g., over North America and large parts of Asia), and an underestimation at lower latitudes during DJF and across most regions during JJA. The overestimation of extremes is linked to advection (primarily in the extratropics) and diabatic processes, with the relative influence of these processes varying across events and regions. This aligns with recent studies indicating an overestimation of daily temperature extremes due to advection and amplified diabatic heating from sensible heat fluxes (Röthlisberger et al., 2025). In contrast, the underestimation in the tropics is primarily associated with biases in the adiabatic contribution and, consequently, in vertical motion, likely stemming from inadequate representation of convective processes and turbulent fluxes (Bao & Stevens, 2021; Bergman & Sardeshmukh, 2004; Stohl, 1998). Similarly, during JJA, the widespread underestimation across most regions results mainly from a combination of underestimated advective and adiabatic processes and also overestimated diabatic processes. These seasonal biases have important implications for targeted model development and improvement.

L245: When mentioning “Rossby wave propagation,” consider providing supporting evidence or omitting the comment.

**Response:** We have revised the wording to avoid any reference to Rossby waves.

L253: I was stumbling across the formulation “This reduction is because ...”

**Response:** We have improved this as “This advective reduction is related to the fact that the future warming at the origin of the traced air masses (at -3d) is 9.4°C at  $t-1$  but only 8.3°C at  $t$  (Figure 5n)”.

L585: “driven by” and “due to” as before. I think you should refrain from implying causality here. Several transition words (e.g., “conversely” in L627, “however” in L583, and “in contrast” in L660) do not seem appropriate in their current context and may be misleading. Revisiting these connectors could help improve clarity and flow.

**Response:** We have revised causal phrasing throughout and reconsidered transition words to improve clarity and accuracy.

## Response to Reviewer 2

We would like to thank the reviewer for their helpful comments. Our responses are printed in blue, while the reviewer's questions are in black. In addition to addressing the individual comments, we have reviewed the manuscript for clarity and flow.

The paper is clear and well-written (although some parts have extensive descriptions of the atmospheric circulation that could be shortened in my opinion). I was already a reviewer in the first part of this article, and I have nothing new to add to the methods, which are essentially similar methods. I have some technical comments below, including some statistical significance computations that should be done differently in my opinion, but apart from those, I would be happy to recommend the paper after some revisions.

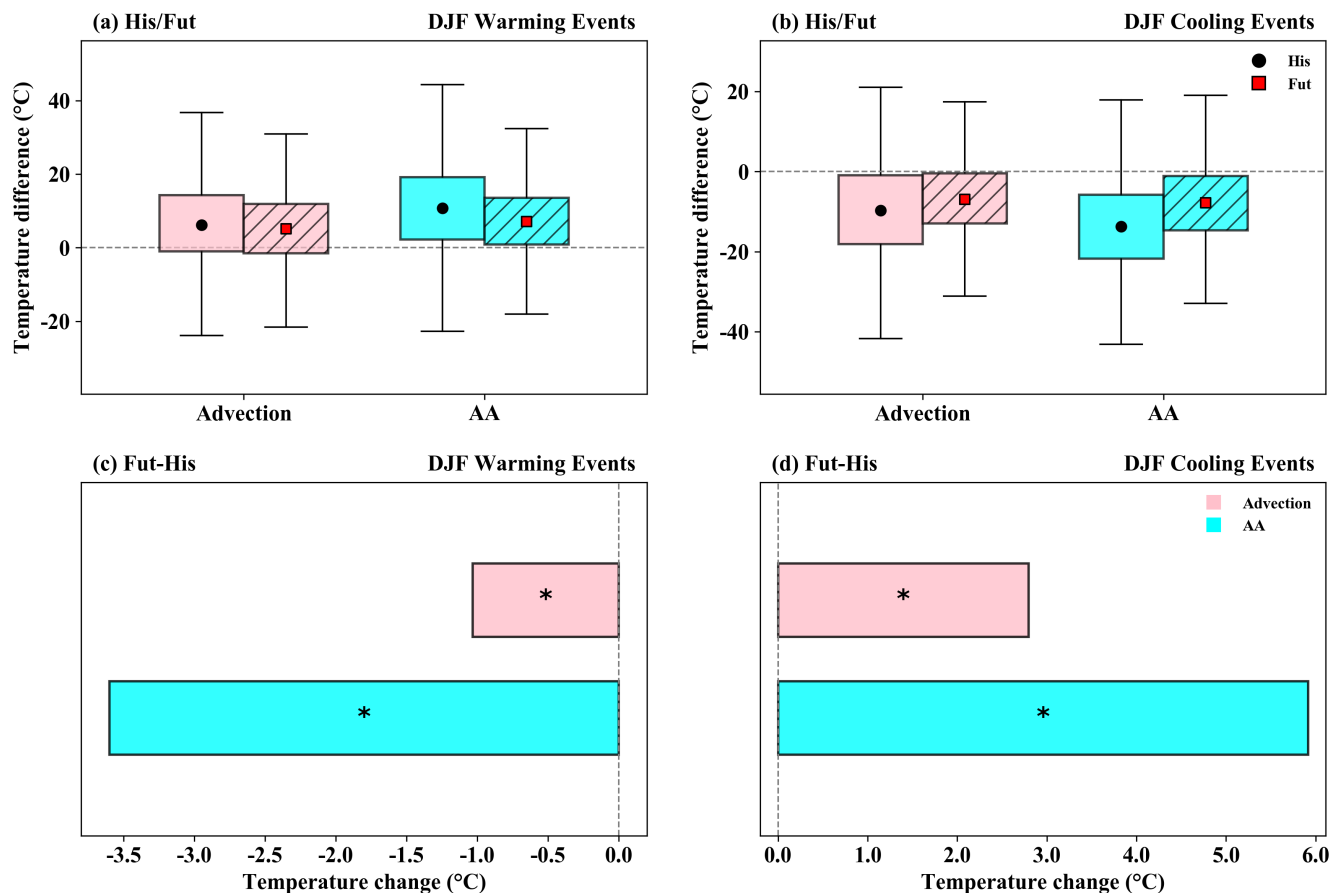
**Response:** Thank you for your continued review and positive assessment of our paper. We appreciate your suggestion to shorten atmospheric circulation descriptions where possible and have streamlined these sections accordingly in the revised version.

Regarding the statistical significance computations, we agreed that this is an important point and have implemented the bootstrapping method with a False Discovery Rate (FDR) test as suggested. The results are shown in Figure R2, and we have revised the manuscript accordingly.

The only main limitation that I see is that I find the paper very descriptive without testing any physical theory. In other words, it could have been interesting to explore how the changes you see fit with some physical expectations for how those mechanisms are supposed to evolve. It is mentioned several times that Arctic amplification, and the associated change in the temperature gradient and general circulation, is expected to decrease the importance of advection, and I think this kind of reasoning could be interesting to investigate further. I do not think this is a reason to reject the paper, but it could really add something on top of those descriptive mechanisms.

**Response:** Thank you for raising this important point regarding the connection between our results and the core concept of Arctic amplification. We have now incorporated a more detailed analysis to address this directly and added these results to the supplementary materials.

Our trajectory analysis reveals a clear fingerprint of Arctic amplification: air originating from the Arctic experiences consistently larger future warming than air from lower latitudes (Figure 5, main paper). This reflects a weakening of the meridional temperature gradient. Furthermore, as shown in Figure R6, air masses on both day  $t-1$  and day  $t$  are substantially modified by the warming climate. The ratio of local temperature change to global warming exceeds 1 on all days ( $-72$  h to 0 h), and is above 1.4 for both day  $t-1$  for warming and day  $t$  for cooling events, satisfying the quantitative definition of Arctic amplification. This confirms that air masses participating in these extreme DTD events carry the signature of Arctic amplification, as they originate from regions warming at a rate significantly faster than the global average.



**Figure R1.** (a,b) Box plots comparing the contribution of advection and the estimated effect of Arctic Amplification (AA) for (a) warming events and (b) cooling events in North America, on a 3-day time scale. For each process, the left box (solid edge) represents the historical climate, and the right box (hatched) represents the future climate. The box spans the 25th to 75th percentiles; the black dot marks the historical mean; the red square marks the future mean; whiskers extend to 1.5 times the interquartile range. The effect of AA is estimated by comparing the seasonal-mean projected temperature change at the origins (at -72 h) of the historical backward trajectories for days  $t$  and  $t-1$ . (c,d) Stacked horizontal bars showing the mean change (future minus historical) for advection and AA for (c) warming and (d) cooling events. The \* indicates that the difference is statistically significant at the 95% confidence level based on a t-test.

We isolated the Arctic Amplification (AA) contribution to extreme temperature events by comparing the seasonal-mean projected temperature change at the origins (at -72 h) of historical backward trajectories from day  $t$  to day  $t-1$  in Figure R1. In this way, only the warming at the origins (which differs between  $t$  and  $t-1$  due to the typically different latitudes of air parcels at -72h) is considered under the assumption that the circulation is unchanged. The magnitudes of the AA effects are more pronounced than those from the changes in advection determined from equation 2. Specifically, for warming events, the change in advection ( $-1.1$  °C) is smaller in magnitude than the change in AA ( $-3.5$  °C), while for cooling events, the change in advection ( $+2.8$  °C) is also smaller than the change in AA ( $+5.9$  °C). The effect of AA explains the weakening of the advective contribution to both warming and cooling events, but the magnitude of the AA effect is mitigated by circulation changes. Also, the event-to-event variability (boxes and whiskers in Figure R1) is very similar between advection and estimated AA contributions.

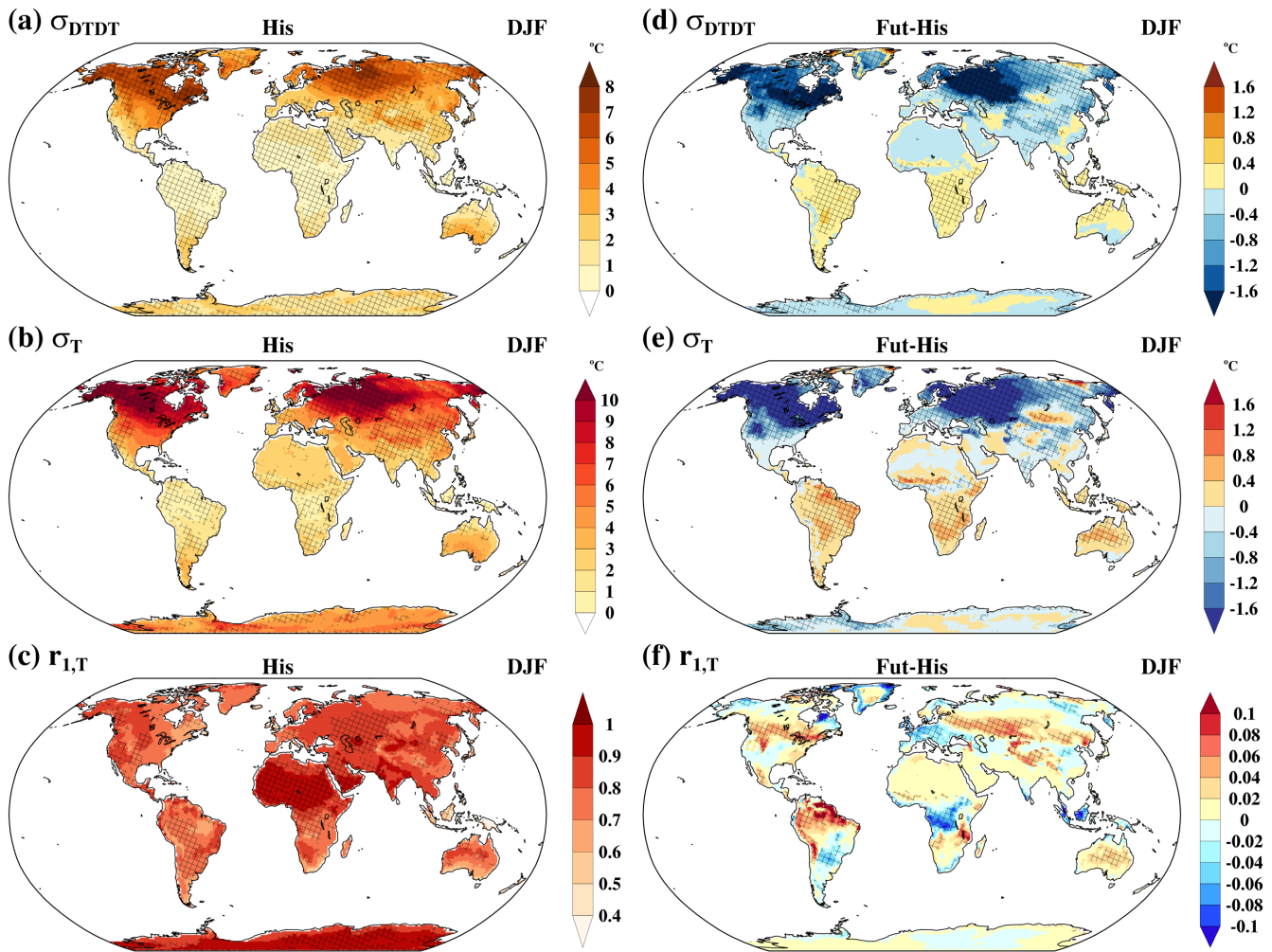
We have added the Figure to the supplementary material and discussed this result in the main paper.

### Major comments

1. A. Figure 1 and corresponding: the way the stippings are computed does not look like a proper statistical test to me. a. If I understood correctly, for the present, the authors flag as “significant” the grid points where 80% of members are within +/- 10% of the ERA5-derived respective quantities. I do not think this is correct: first the +/-10% for ERA5 is an ad-hoc measure of the uncertainty. Second, I do not see why 80% of the members should be a correct measure of a significant difference. I suggest to do an actual statistical test with a standard reference level of 95% significance for example. In essence you want to know whether the climate of CESM is compatible with the value for ERA5: that is what you need to test for. The climate of CESM is defined by all the members being put together: what you need to test is whether the  $\sigma\_DTDT$ ,  $\sigma\_T$  and  $r\_1,T$  of this climate are compatible with the same values for ERA5 (which also has an uncertainty). You could for example, employ a bootstrapping approach on the ERA5 data (the values for the model are likely very well estimated given the amount of members you have) and check whether the distribution of values you obtain are compatible with the one from the model at the 95% level.

**Response:** Thank you for your valuable feedback on our significance testing methodology. Following your suggestion, we have revised our approach (Figure R2). Given the present-day climate, we use bootstrapping ( $n=1000$ ) on ERA5 data to construct an observational uncertainty distribution, which we compare with the CESM multi-member distribution at the 95% significance level. As noted in Comment 3, we apply a false discovery rate (FDR) correction to the bootstrapping results. We found that this FDR correction reduces the significant area by only 1–8% globally at the 1% level (Figure R2a-c). We have updated the methodology and results accordingly.

We have explained this result in the revised manuscript as follows: “The hatching in Figure R2a-c indicates that the CESM ensemble mean is systematically different from the ERA5 estimates in most regions, with only a few exceptions. Across approximately 40–70% of global land areas, CESM-LE exhibits systematic biases in daily temperatures and related metrics relative to ERA5. While the model reproduces the large-scale spatial patterns of the metrics (see Part I), notable discrepancies remain in the magnitude of anomalies (CESM-ERA5, Figure R3). These biases include an overestimation of DTDT extremes, particularly in the mid-to-high latitudes such as over North America and large parts of Asia. In contrast, underestimations occur at lower latitudes (e.g., over the Sahara Desert and Amazon) during DJF and persist across most regions during JJA.



**Figure R2.** The ensemble means of (a, d) standard deviation of DTD variations ( $\sigma_{DTDT}$ , °C), (b, e) standard deviation of daily mean temperature ( $\sigma_T$ , °C), and (c, f) lag-1 autocorrelation of daily mean temperature ( $r_{1,T}$ ) in December–February (DJF) in the historical climate (a–c) and projected change (d–f). In the panels a–c, cross-hatching marks grid points where the CESM-LE ensemble mean differs significantly from the ERA5-derived metric, with statistical significance determined through bootstrap resampling. In panels d–f, cross-hatching denotes grid points where the future minus historical difference is significantly different from zero, assessed via a two-sample bootstrap test.

1. B. Same for the future: why don't you simply test whether there is a significant difference in  $\sigma_{DTDT}$ ,  $\sigma_T$  and  $r_{1,T}$  between the two climates by putting all members together in each period?

**Response:** Thank you for this helpful suggestion. For the future projections, we have adopted the same approach: test whether there is a significant difference in temperature metrics between the historical and future climates by pooling all ensemble members from each period (Figure R2d-f). We have applied bootstrapping with FDR correction at the 1% and have updated the result accordingly.

2. For all your significance maps, you need to take into account correlations in statistical testing and employ a false discovery rate, see Wilks (2016).

**Response:** Thank you for your suggestion. Following your suggestion, we applied the bootstrapping test with false discovery rate (FDR) correction (Wilks, 2016) to all climatological maps, while for the other composite maps, we applied a t-test at the 95% significance level. After applying the FDR correction, the area with significant differences is reduced only by 1–8%.

3. Several times, the authors argue that the model is doing a reasonable job in reproducing the statistics of ERA5. I am not sure this is so much the case, as exemplified by Figure 1, for example, where the stipplings do not really cover most of the regions (modulo my main comment 1). I think you should emphasise the differences more, including quantifying them when possible. One point, for example, is that the model seems to have a diabatic contribution larger than ERA5, which is something also found recently by Röthlisberger et al. (2025) in a different context: it seems to me that the model may be right for the wrong reasons.

**Response:** Thank you for the suggestion. We agree that our initial framing was overly generous and have revised the manuscript to place greater emphasis on the quantitative differences between the model and ERA5. We have compared our results with those of previous studies by Röthlisberger et al. (2025) and incorporated recent findings into the Discussion.

Approximately 40–70% of global land regions exhibit systematic biases in daily temperature and related metrics simulated by CESM-LE compared with ERA5 (Figures R2–R3). While the model captures large-scale spatial patterns (cf. Part I), there are notable discrepancies in the magnitude of anomalies (Figure R3) and in the contributions of underlying physical processes. These biases include an overestimation of DTD extremes, most pronounced in the mid-to-high latitudes (e.g., over North America and large parts of Asia), and an underestimation at lower latitudes during DJF and across most regions during JJA. The overestimation of extremes is linked to advection (primarily in the extratropics) and diabatic processes, with the relative influence of these processes varying across events and regions. This aligns with recent studies indicating an overestimation of daily temperature extremes due to advection and amplified diabatic heating from sensible heat fluxes (Röthlisberger et al., 2025). In contrast, the underestimation in the tropics is primarily associated with biases in the adiabatic contribution and, consequently, in vertical motion, likely stemming from inadequate representation of convective processes and turbulent fluxes in the models (Bao & Stevens, 2021; Bergman & Sardeshmukh, 2004; Stohl, 1998). Similarly, during JJA, the widespread underestimation across most regions results mainly from a combination of underestimated advective and adiabatic processes and also overestimated diabatic processes. These seasonal biases have important implications for targeted model development and improvement.

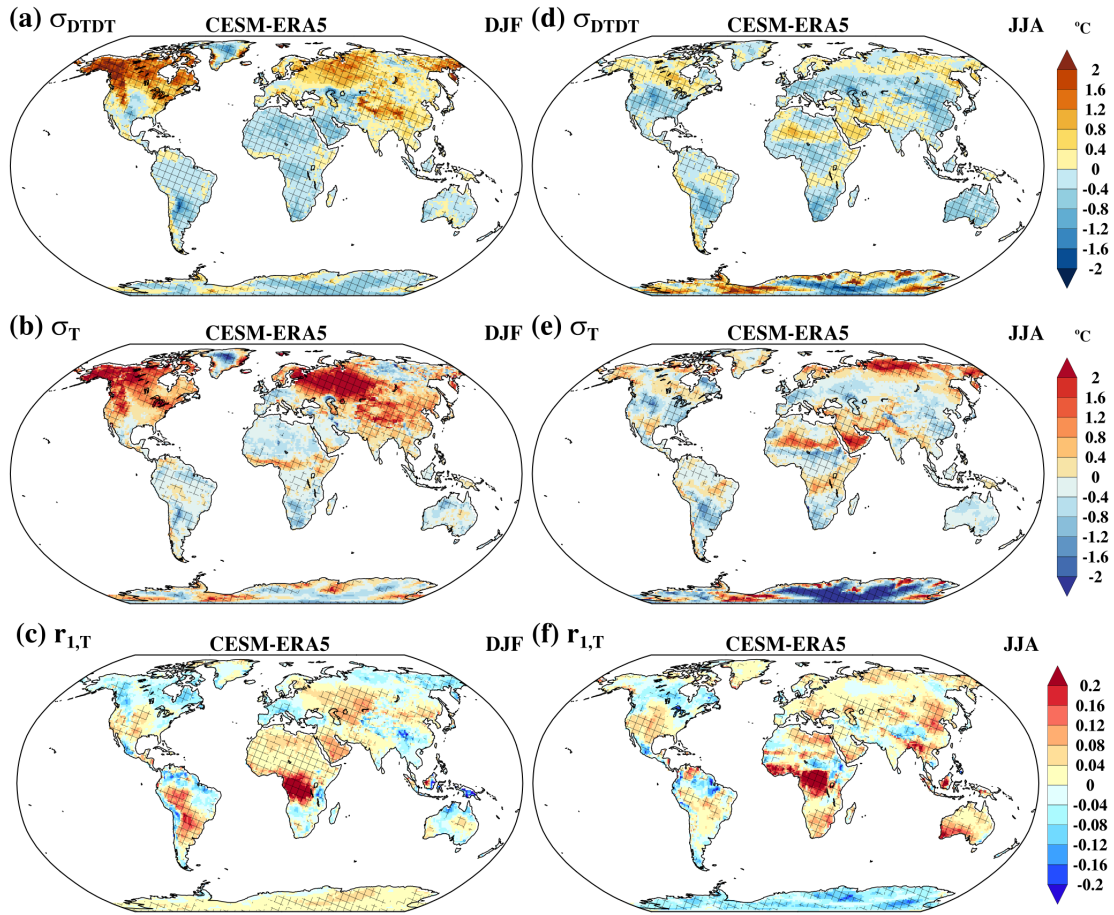


Figure R3. The absolute differences between the CESM-LE and ERA5 of (a, d) standard deviation of DTTD variations ( $\sigma_{DTDT}$ , °C), (b, e) standard deviation of daily mean temperature ( $\sigma_T$ , °C), and (c, f) lag-1 autocorrelation of daily mean temperature ( $r_{1,T}$ ) in December-February (DJF, a-c) and June-August (JJA, d-f). In the panels a–d, cross-hatching marks grid points where the CESM-LE ensemble mean differs significantly from the ERA5-derived metric, with statistical significance determined through bootstrap resampling.

#### Minor comments

1. Please precise which version of CESM you are using.

**Response:** We have used 1st version of CESM-LE and have corrected this in the revised manuscript.

2. Figure 1: Because of the strong meridional differences, you could plot the changes in the second column in percentage rather than absolute values.

**Response:** Thank you for the suggestion. We have plotted the projected percentage change in Figure R4. However, to maintain consistency with the trajectory analysis for which absolute changes are much easier to interpret, we have kept the absolute differences in the main paper and added this Figure to the supplement.

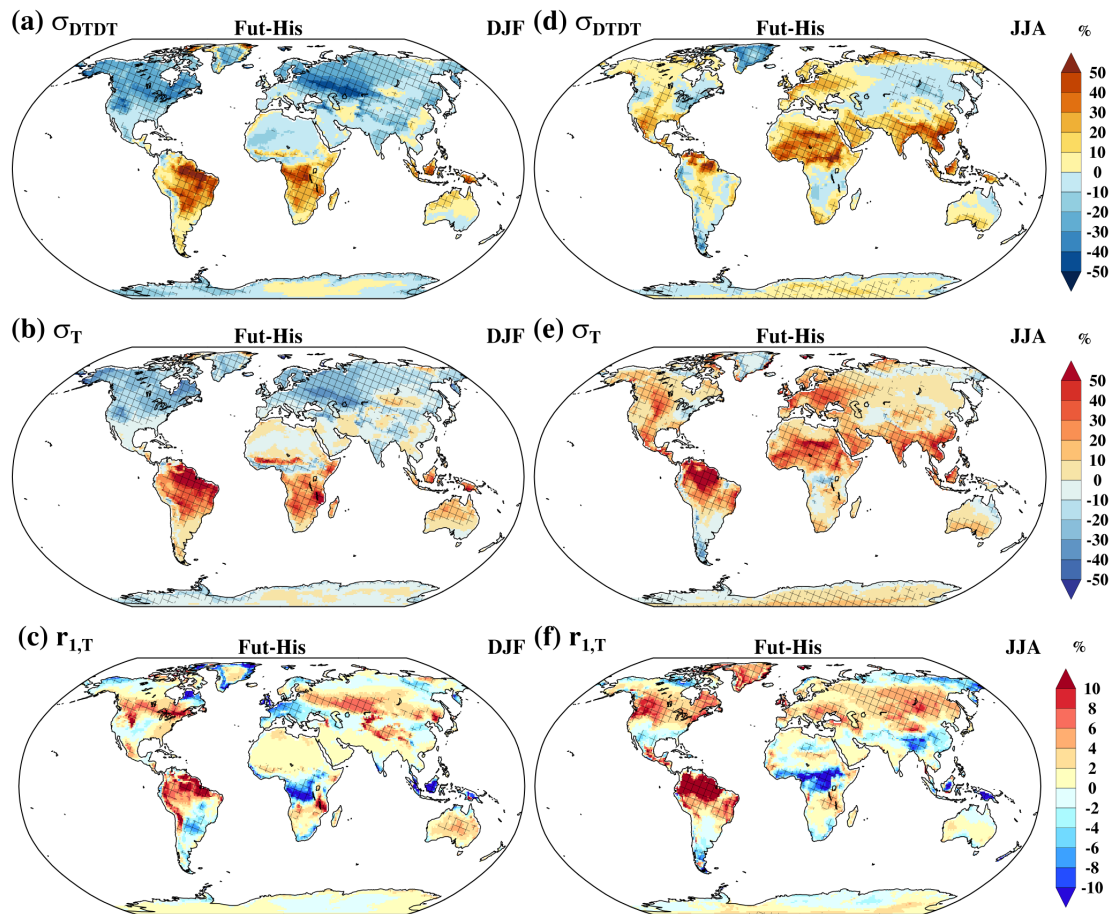


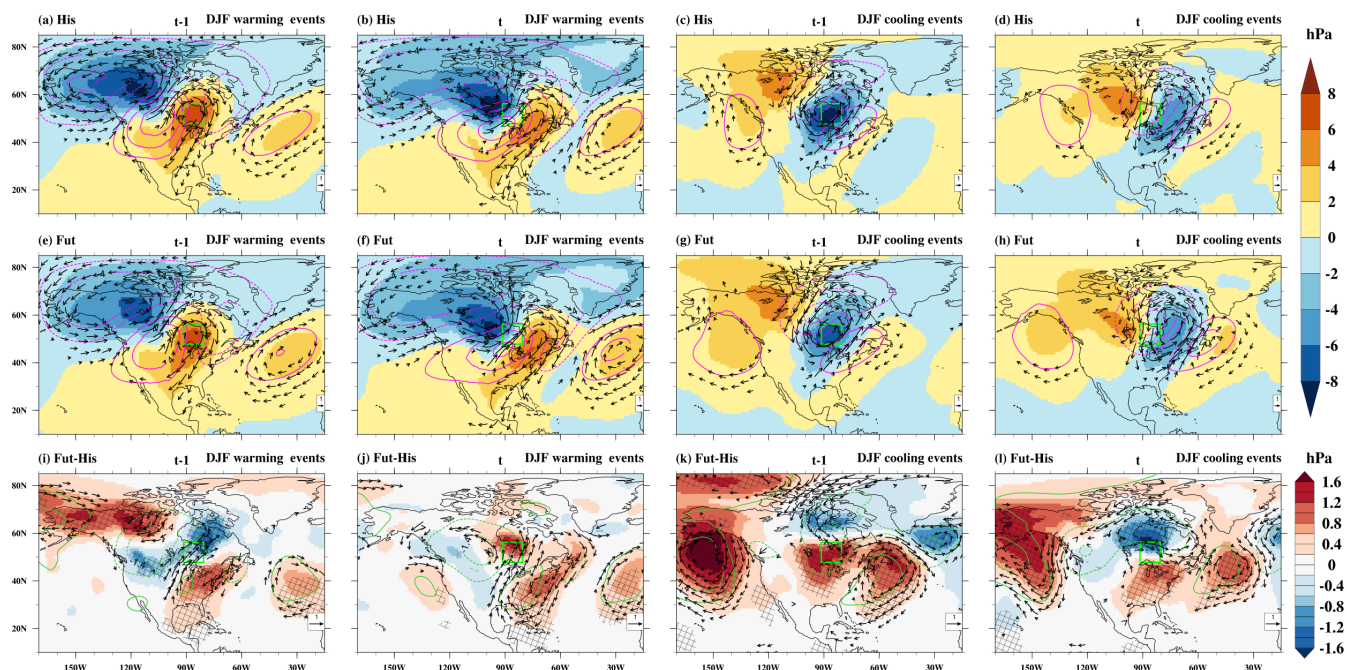
Figure R4. The projected percentage changes of (a, d) standard deviation of DTTD variations ( $\sigma_{DTDT}$ , %), (b, e) standard deviation of daily mean temperature ( $\sigma_T$ , %), and (c, f) lag-1 autocorrelation of daily mean temperature ( $r_{1,T}$ , %) in December-February (DJF, a-c) and June-August (JJA, d-f). In panels a–f, cross-hatching denotes grid points where the future-minus historical difference is significantly different from zero, assessed via a two-sample bootstrap test.

3. Figure 3: The sigma\_DTTD should be delta\_T?

**Response:** We have changed this.

4. Figure 4: The stiplings are barely visible.

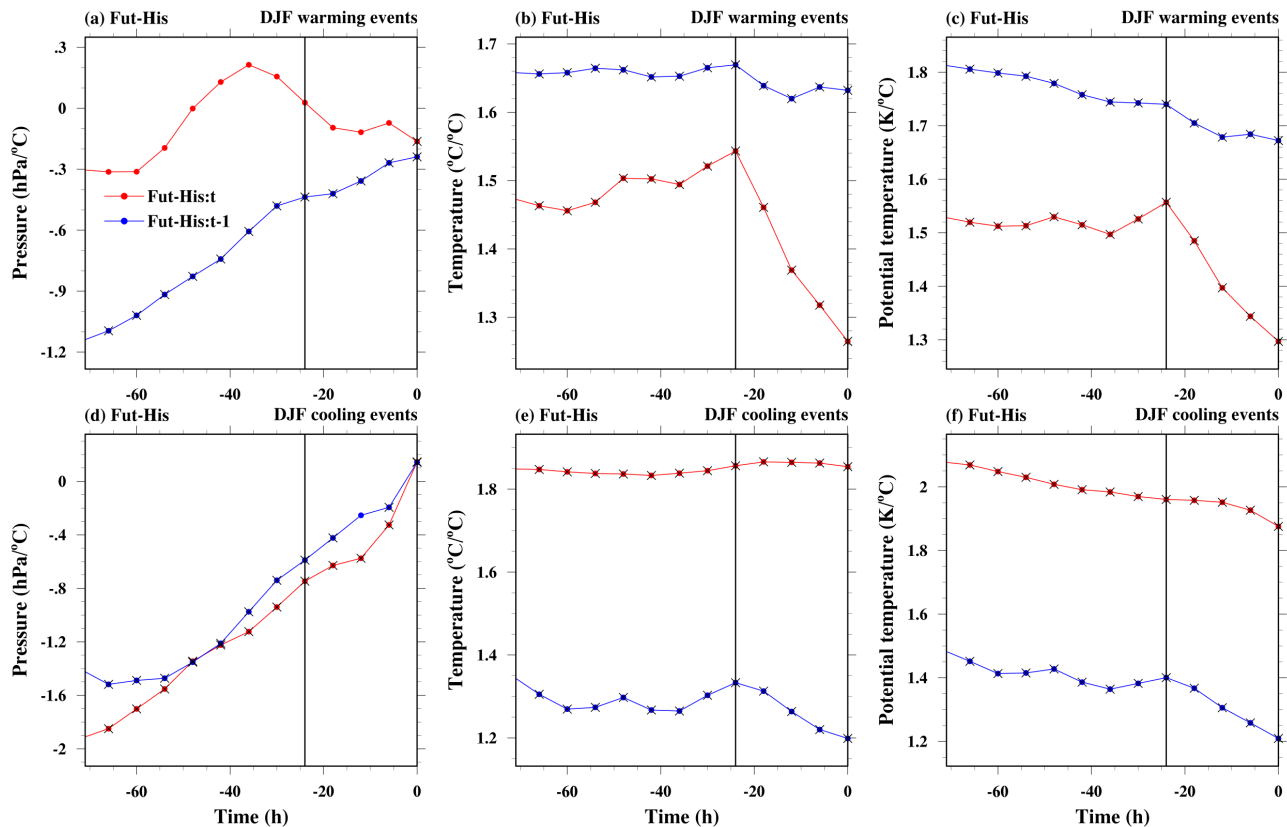
**Response:** Thank you for the suggestion. We have now changed the stipling to cross-hatching, which is more visible, and reduced the spacing between the panels, as shown in Figure R5 below.



**Figure R5.** Composite of sea level pressure anomalies (hPa, color shading), wind anomalies at 850 hPa ( $\text{m s}^{-1}$ , vectors), and geopotential height anomalies at 500 hPa (gpm, magenta and darkgreen contours) relative to the seasonal mean on the (a, e, i, c, g, k) previous day ( $t-1$ ) and (b, f, j, d, h, l) the event day ( $t$ ) of the warming (a-b, e-f and i-j) and cooling (c-d, g-h and k-l) events during December-February (DJF) in (a-d) historical climate (His), (e-h) future climate (Fut), and (i-l) projected future changes (Fut-His) at a selected grid box in North America (green box). Note that, in (a-h), wind vector anomalies  $\geq 2 \text{ m s}^{-1}$  and in (i-l), wind vector difference anomalies  $\geq 0.5 \text{ m s}^{-1}$  are plotted. The dotted and bold contours indicate negative and positive geopotential height anomalies, respectively. Additionally, the cross-hatching area indicates where the ensemble mean of sea level pressure differences exceeds the 95% confidence threshold based on a t-test.

5. Figure 5: I would suggest scaling the temperature and pressure differences by a global/regional warming level to see what is changing beyond the expected local warming.

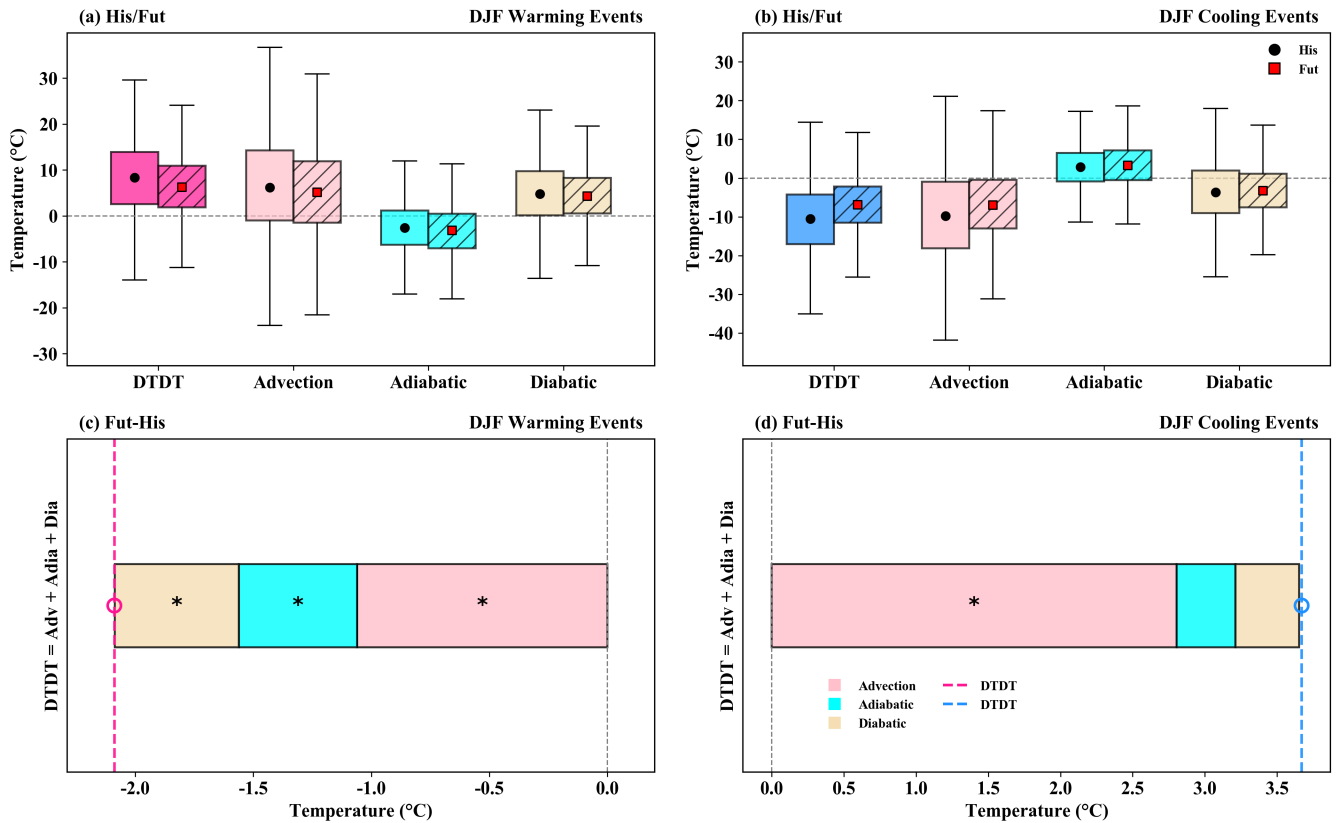
**Response:** Thank you for the excellent suggestion. We have scaled the projected changes in temperature, pressure, and potential temperature with annual global temperature change, as shown in Figure R6. This updated analysis also illustrates how circulation patterns and temperatures are influenced by warming. The ratio of local temperature change to global warming exceeds 1 on all days ( $-72 \text{ h}$  to  $0 \text{ h}$ ) and is above 1.4 for both day  $t-1$  warming and day  $t$  cooling events, satisfying the quantitative definition of Arctic amplification. This confirms that air masses participating in these extreme DTD events carry the signature of the Arctic amplification signal, as they originate from regions warming at a rate significantly faster than the global average.



**Figure R6:** The mean Lagrangian evolution of distinct physical parameters: (a, d) pressure (hPa), (b, e) temperature ( $^{\circ}\text{C}$ ), and (c, f) potential temperature (K) is shown along the air mass trajectories initialised on the previous ( $t-1$ ) and event ( $t$ ) days for projected changes in extremes scaled by annual global surface temperature ( $^{\circ}\text{C}$ ). Additionally, Bold circles with crosses mark time steps where the projected changes are statistically significant at the 95 % confidence level based on a t-test.

6. Figure 6 and similar: maybe you could add the boxplots for the future on panels a and b, also to compare the spread in each period (I do not expect the spread to be small, thus the changes you observe are probably much smaller in intensity compared to the spread between events in each period).

**Response:** We thank the reviewer for this constructive suggestion. We have updated Figure 6 (and other relevant figures) to include boxplots for both the historical and future periods in panels a and b. The spread among events within each period is indeed substantial. For the grid point over North America shown in Figure R7, the future spread is slightly smaller than (for DTD, advection, and diabatic heating) or similar to (for adiabatic processes) the event-to-event variability in the historical climate, suggesting that this spread may also change in a warmer climate.



**Figure R7.** The contribution of the different physical processes (advection, adiabatic and diabatic temperature change) over North America during December-February (DJF) to genesis of DTDT (a, c) warming and (b, d) cooling events during historical/future climate (a-b, box plots) and projected future change (c-d, stacked plots) according to Eq. (2), which refers to a 3d-time scale. The box spans the 25th and 75th percentiles of the trajectory data; the black dot inside the box gives the mean of the related quantities in the historical climate, and the whiskers indicate 1.5 times the interquartile range in panels (a) and (b). The dotted lines in the stacked plots in panels (c) and (d) show the mean future change for DTDT warming and cooling events, respectively, and coloured bars indicate the contributions of the individual processes. Circle and \* symbols mark future change distributions for which the ensemble mean differences exceed the 95% confidence threshold based on a t-test.

7. Figure 7: Why did you decide to change the position of the box for looking at extreme DTDT changes compared to Figure 4?

Response: The grid points were chosen based on regions exhibiting the most significant changes (Figure R4). For DJF, the largest variations occurred over North America and Western North America, and the mechanisms behind the extremes are similar. For JJA, a grid point in western North America was selected, as the largest significant decrease was located near the coast (Figure R4d) and followed a different mechanism than at the northern grid point.

## References

- Bao, J., & Stevens, B. (2021). The Elements of the Thermodynamic Structure of the Tropical Atmosphere. *Journal of the Meteorological Society of Japan. Ser. II*, 99(6), 1483-1499. <https://doi.org/10.2151/jmsj.2021-072>
- Bergman, J. W., & Sardeshmukh, P. D. (2004). Dynamic Stabilization of Atmospheric Single Column Models. *Journal of Climate*, 17(5), 1004-1021. [https://doi.org/https://doi.org/10.1175/1520-0442\(2004\)017](https://doi.org/https://doi.org/10.1175/1520-0442(2004)017)
- Bieli, M., Pfahl, S., & Wernli, H. (2015). A Lagrangian investigation of hot and cold temperature extremes in Europe. *Quarterly Journal of the Royal Meteorological Society*, 141(686), 98-108. <https://doi.org/https://doi.org/10.1002/qj.2339>
- Hamal, K., & Pfahl, S. (2025). Physical processes leading to extreme day-to-day temperature change – Part 1: Present-day climate. *Weather Clim. Dynam.*, 6(3), 879-899. <https://doi.org/10.5194/wcd-6-879-2025>
- Röthlisberger, M., & Papritz, L. (2023). Quantifying the physical processes leading to atmospheric hot extremes at a global scale. *Nature Geoscience*, 16(3), 210-216. <https://doi.org/10.1038/s41561-023-01126-1>
- Röthlisberger, M., Sprenger, M., Beyerle, U., Fischer, E. M., & Wernli, H. (2025). Advective, adiabatic and diabatic contributions to heat extremes simulated with the Community Earth System Model version 2. *EGUsphere*, 2025, 1-32. <https://doi.org/10.5194/egusphere-2025-5146>
- Stohl, A. (1998). Computation, accuracy and applications of trajectories—A review and bibliography. *Atmospheric Environment*, 32(6), 947-966. [https://doi.org/https://doi.org/10.1016/S1352-2310\(97\)00457-3](https://doi.org/https://doi.org/10.1016/S1352-2310(97)00457-3)
- Wilks, D. S. (2016). “The Stippling Shows Statistically Significant Grid Points”: How Research Results are Routinely Overstated and Overinterpreted, and What to Do about It. *Bulletin of the American Meteorological Society*, 97(12), 2263-2273. <https://doi.org/https://doi.org/10.1175/BAMS-D-15-00267.1>