

1 Dear Editor,

2 Thank you for your valuable feedback. We have carefully considered both of the major comments raised in the  
3 revision process and addressed all of the reviewers' questions in detail in the updated version of the manuscript.

#### 4 **0.1 General comments:**

5 **A major revision is needed to improve the scientific quality and presentation quality of this manuscript.**  
6 **As pointed out by one of the reviewers, the conceptual clarity of the evaluation approach and the**  
7 **linkages between different parts of the study need to be addressed. MAE reductions can be caused**  
8 **by different reasons, such as by reductions in the climatological mean bias, by reductions in the bias**  
9 **of the variance, by improvements of correlations between simulated and observed time series, or by**  
10 **a combination of these factors. The authors should consider including multiple metrics to evaluate**  
11 **different aspects of the model and improving the coherency of the manuscript.**

12 To support the MAE reductions reported in our CMIP6-CORDEX-CA simulations, we have incorporated ad-  
13 ditional analyses. Specifically, we now evaluate the model's performance using multiple metrics, including MAE,  
14 RMSE, climatological bias, correlation, comparison of probability density functions (PDFs) of maximum daily pre-  
15 cipitation (provided in the Appendix), and the Relative Probability Difference Index. These additions provide a  
16 more comprehensive assessment of the error reduction mechanisms.

17 Furthermore, we conducted a comparative analysis of the two GCM simulations used in our study to explore  
18 their similarities.

19 We discussed why the unconstrained CNN might have performed better.

20 In the revised manuscript, we have also provided an extended discussion on the caveats, strengths, and future  
21 outlook of the work.

22 On behalf of all the authors, Bijan Fallah

1 Dear Editor and Reviewers,

2 Many thanks for your valuable comments and suggestions, which were very useful for improving presentation of  
3 results and paper readability. In the following, we will answer all the comments raised by the reviewers in detail.  
4 The reviewers' comments are in bold, citations in *Italic* and our answers are in regular font.

## 5 **1 Reviewer 1:**

### 6 **1.1 Minor comments**

7 **Good job for the improvement in the writing quality of your manuscript; it's now a lot better! I**  
8 **spotted some last mistakes and missing words that you'll find in my annotated manuscript with**  
9 **comments (see annotated pdf). The other comments in the annotated pdf apply to results and**  
10 **discussion, which need some more work to present your results fairly.**

11 Thank you very much for your valuable comments. We will first address the comments provided in the PDF file,  
12 followed by a response to the remaining feedback:

13 **P1: be careful with using the terms significant as it requires a statistical test of significance. If**  
14 **you didn't do that I would use another word that is a synonym.**

15 Done.

16 **P1: it's not clear here what this value represents and I also cannot find it in the results. Is it a**  
17 **difference between in MAE between the RCM and GCM or something else? This should be specified.**

18 That is indeed very true. Thanks for hinting that. We changed the text to clarify the reduction in MAE and  
19 Bias during summer, winter and annual. Also we added some analysis about the frequency of maximum daily values  
20 to the Appendix A. That analysis is also linked to the editor's comment. We show that higher resolution datasets  
21 show maximum daily values nearer to the CHIRPS dataset and the Global models or reanalysis are not able to do  
22 so.

23 **P1: see comment above, and provide a number.**

24 We changed the text accordingly. We added the following information to the abstract : "The number of days  
25 with precipitation exceeding 20 mm increases by more than 90 by the end of the century, compared to the historical  
26 reference period, under the SSP3-7.0 and SSP5-8.5 scenarios. The annual 99th percentile of total precipitation  
27 increases by more than 9 mm/day over mountainous areas of Central Asia by the end of the century, relative to the

28 1985–2014 reference period, under the SSP3-7.0 and SSP5-8.5 scenarios.”

29 **P1: mmh you showed that it actually has trouble generalising to a new GCM so how can it have**  
30 **an added value?**

31 Done! The new text is : ”The CNN successfully emulates the GCM-CCLM model chain over large areas of CA,  
32 but shows reduced skill when applied to a different GCM-CCLM model chain.”

33 **P3: that is implied**

34 We removed this sentence.

35 **P6: change this to ”imperfect model” setup. The ”this” is not clear to which it refers**

36 We used ”In the imperfect model setup,”

37 **P6: the Done!**

38 **P6: did you use cross validation?**

39 No, CR is not commonly used in deep learning approaches (like CNNs) because it is computationally expensive.  
40 Instead, we rely on a single validation set to track model performance. This is due to the large number of parameters  
41 in CNNs and the resources required to train these models over multiple folds (Bengio, 2012).

42 While cross-validation can help with smaller datasets, there is also a risk of overfitting if the model is highly  
43 complex (like CNNs). Since CNNs are powerful enough to memorize smaller datasets, using the same data in  
44 multiple folds could lead to overfitting as the model is exposed to repeated patterns. In such cases, cross-validation  
45 might not prevent overfitting but could actually increase its likelihood.

46 In our case, if the dataset is small, the risk of overfitting is more related to the capacity of the CNN rather than  
47 the absence of cross-validation. Cross-validation helps small datasets but can overfit if the model complexity is too  
48 high compared to the size of the data (Goodfellow, 2016).

49 To demonstrate that the out training step does not trap in this problem, we show the training and validation  
50 loss diagram in Figure 1. As can be seen the validation and training loss reaches a steady state at the end of the  
51 training process for NoneCL (other models show similar behavior).

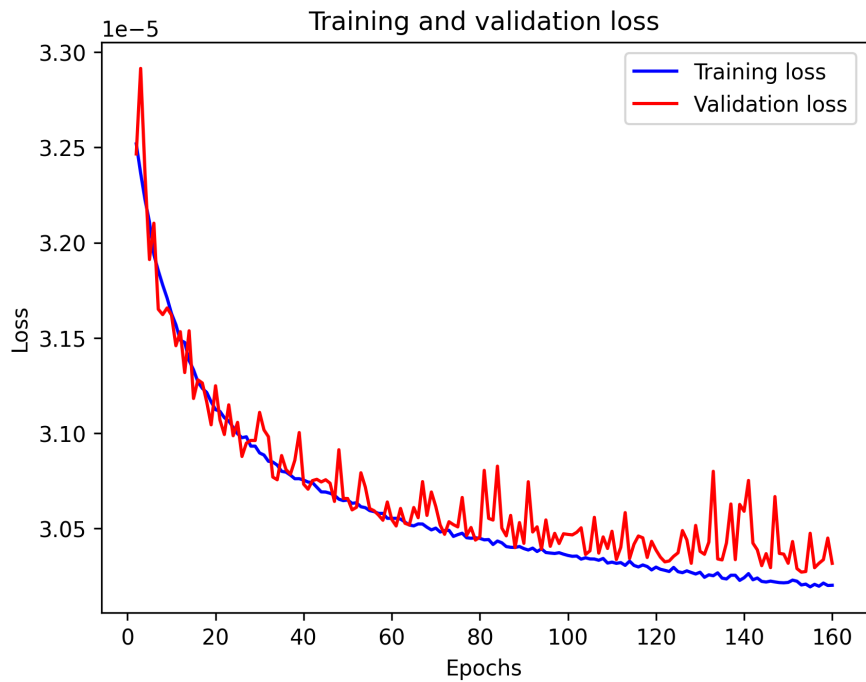


Figure 1: Training and validation loss with respect to epochs.

52 Overall, this seems like a reasonable training curve for a CNN. The model shows steady improvement on both  
 53 training and validation data without over-fitting.

54 **P6: independent**

55 Done!

56 **P6: this should follow the order in which they come below**

57 We reordered them.

58 **P7: function**

59 Done!

60 **P7: "Then the mean absolute error (MAE) is applied as a loss function between the final output  
 61 and the target value."**

62 Corrected!

63 **P7: It's not very clear here if there is a constraint layer like hard constraining or just a special  
 64 loss.**

65 This approach does not impose a hard constraint through a dedicated layer or architectural component. Instead,

66 we add a regularization term to the loss function, which allows for soft enforcement of the constraint. The Mean  
67 Absolute Error (MAE) loss is extended with an additional constraint violation (CV) loss term.

68 We clarified that by adding those information in the new version of the manuscript. Thanks for this comment.

69 **P7: The model is trained over**

70 added.

71 **P7 : Can you give a reason why the lr for SCL is so much smaller? There's a factor 100 between**  
72 **them.**

73 It is purely based on the empirical performance (we have tested several  $lr$  values). A lower learning rate allows the  
74 model to fine-tune more carefully in the presence of complex loss terms, leading to better convergence and improved  
75 results. This is likely why we observed a better performance with this setup.

76 Harder et al. (2023) also used different learning rates for their constrained and unconstrained models. While they  
77 did not explicitly state a reason for using a much lower learning rate for the soft-constrained loss (SCL), this is a  
78 common strategy in deep learning to stabilize training in cases where the additional loss term (such as a regularization  
79 or constraint violation term) can lead to more significant gradients or slower convergence.

80 **P7: This is superfluous, especially because evaluation comes later. So you can remove this.**

81 This part was added due to your comment in the previous round (on the usage of MAE). We have removed it  
82 from the new version now.

83 **P8: Above you used the term "validation", stay consistent on the terms.**

84 Done!

85 **P8: this sentence is unclear, using CHIRPS directly as what? Are these the 20% of testing stated**  
86 **above in l. 164? Needs rewriting and I think you can remove the "instead of [...]" directly**

87 We have modified the paragraph and splitted that into two paragraphs: *"According to Ciarlo et al. (2021),*  
88 *the choice of observational data can significantly influence the perceived added value of an RCM, particularly when*  
89 *detecting extreme events, where poor-quality data might misleadingly suggest improved model performance. They*  
90 *recommend using observational datasets with spatiotemporal resolutions comparable to the model's for enhanced ac-*  
91 *curacy. In line with this, we use CHIRPS, a high-resolution gridded observational dataset, to validate the CCLM*  
92 *driven by the GCM. CHIRPS offers a resolution of  $0.05^\circ$ , covering latitudes from  $50^\circ S$  to  $50^\circ N$ , and provides inde-*  
93 *pendent observations derived from satellite and station data. This contrasts with reanalysis datasets, which rely on*  
94 *climate model simulations (Funk et al., 2015). For the validation of the CNN, however, we allocate 20% of the CCLM*

95 *simulation data as the target for evaluating the CNN emulator's performance rather than directly using CHIRPS.*  
96 *This is because the CNN is designed to emulate the climate output produced by the CCLM, not to match observational*  
97 *data directly. While CHIRPS is used to validate the accuracy of the CCLM output, we validate the CNN by ensuring*  
98 *it accurately reproduces the CCLM's fine-scale climate information, which has already been verified against CHIRPS*  
99 *for its realism."*

100 **P8: Wouldn't F be the same as yi?**

101 We modified that.

102 **P8: as an evaluation metric**

103 Done!

104 **P8: This should be on a new line.**

105 Done!

106 **P8: and this should be F again?**

107 Done!

108 **P8: above this is called O**

109 Done!

110 **P9: can we have a value here? And you are plotting the difference so you don't know if the**  
111 **CCLM's MAE is "high" or "lower", only that it's difference is big.**

112 We are showing the differences of CCLM's MAE and the one from ERAInterim. If we have positive values it  
113 means that errors of ERAInterim are larger and 5 mm/day is also a significant difference. And magenta regions  
114 show negative values which show better skills for ERAInterim. ERAInterim for example during DJF in figure 4.b  
115 has already 2-3 mm/day MAE.

116 We put numbers to clarify that in more details.

117 **what do you mean with increase and reduction of the MAE? Compared to what? I think you**  
118 **should compare it to the GCM. So in Afgh-Taj, the CCLM is closer to observations than the GCM**  
119 **blabla.**

120 We agree with your comment. We will re-frame the text accordingly to : *"In the mountainous areas of*  
121 *Afghanistan, Kyrgyzstan, and Tajikistan, CCLM is closer to observations than GCM. However, GCM is closer*  
122 *to observations near the domain's southern boundaries throughout the year and in the south and southeast during*  
123 *summer"*. In this case GCM is the ERAInterim reanalysis.

124 Other than that, to demonstrate the skill of our CCLM simulations and also in line with the **Editor’s** comment,  
 125 we conducted the following extra analysis, which might be useful for CORDEX-like model evaluation, especially  
 126 those focusing on extreme events like in precipitation. Interestingly, the same added value patterns are obtained by  
 127 the new applied metric, showing our RCM simulation, shows skills similar to a typical CORDEX-CA simulation and  
 128 ads values over mountains of CA.

129 We ”hope” that doing those analysis will show our simulations might be useful for the community:

130 We utilize the spatially distributed added value index as introduced in Ciarlo et al. (2021) to evaluate the  
 131 performance of our RCM simulations. This metric is particularly well-suited as it captures the spatial variability of  
 132 added value, providing a detailed assessment of how well our model reproduces fine-scale climate features compared  
 133 to a coarser global climate model (GCM). By focusing on high-resolution regional models, the added value index  
 134 allows us to quantify improvements in regions with complex topography or localized climate phenomena, which  
 135 are often missed by coarser models. Moreover, it is an ideal metric for examining extreme events, a key focus in  
 136 climate impact studies, as it helps highlight where our emulator outperforms the GCM in simulating such events.  
 137 The spatial nature of this metric also enables us to produce clear visualizations of model performance across various  
 138 regions, facilitating comparison and communication of results. Overall, the added value index offers a robust and  
 139 versatile tool for assessing our model’s ability to capture regional climate details, which is critical for evaluating the  
 140 effectiveness of climate model downscaling.

### 141 **Metrics**

142 We use the formula presented by Ciarlo et al. (2021) to calculate the added value. In the first step, the selected  
 143 GCM, RCM and observational data are interpolated onto the RCM grid using the distance-weighted average method.  
 144 Interpolation of the coarser grid to a higher one might create unrealistic values. This issue was discussed in the work  
 145 of Ciarlo et al. (2021). Then, for each grid point of the domain, we calculate the probability density function (PDF)  
 146 from GCM, RCM and observation. As in the case of Ciarlo et al. (2021), for a fair comparison, the bin size is fixed  
 147 to 1mm/day and the maximum value at each grid point for the calculation of the PDF is taken from the maximum  
 148 of all datasets at that grid point. At each grid point, the absolute differences ( $D$ ) between the frequency ( $N$ ) of the  
 149 model ( $M$ ) and observation ( $O$ ) at each bin ( $\nu_t$ ) are divided by the sum of the observation value:

$$D_M = \frac{\sum_{\nu=1}^{\nu_t} |(N_M - N_O)\Delta\nu|}{\sum_{\nu=1}^{\nu_t} (N_O\Delta\nu)}. \quad (1)$$

150 The added value at each grid point ( $A_i$ ) is then the difference between the  $D_{GCM}$  and  $D_{RCM}$ . Positive (negative)  
 151 values indicate an improvement (degradation) of the downscaling compared to the GCM:

$$A_i = D_{GCM} - D_{RCM}. \quad (2)$$

152 There shall be a condition added to the added value calculation in the case that  $N_{GCM} = 0$  in a bin and  $N_{RCM}$   
 153 and  $N_O$  are nonzero. Under this condition, that bin must not contribute to the  $D_{RCM}$  calculation.

154 Analogous to the  $A_i$  calculation, the climate change downscaling signal is presented by :

$$D_{Mf} = \frac{\sum_{\nu=1}^{\nu_t} |(N_{Mf} - N_{Mh})\Delta\nu|}{\sum_{\nu=1}^{\nu_t} (N_{Mh}\Delta\nu)}. \quad (3)$$

155 where  $Mf$  is the future projection,  $Mh$  the corresponding historical period and :

$$A_{DS} = D_{GCMf} - D_{RCMf}. \quad (4)$$

156 Where  $A_{DS}$  is the downscaling signal, large positive or negative values indicate a significant climate change  
 157 downscaling signal. Values near zero indicate a weak downscaling signal.

#### 158 **Added value of CCLM driven by ERAInterim**

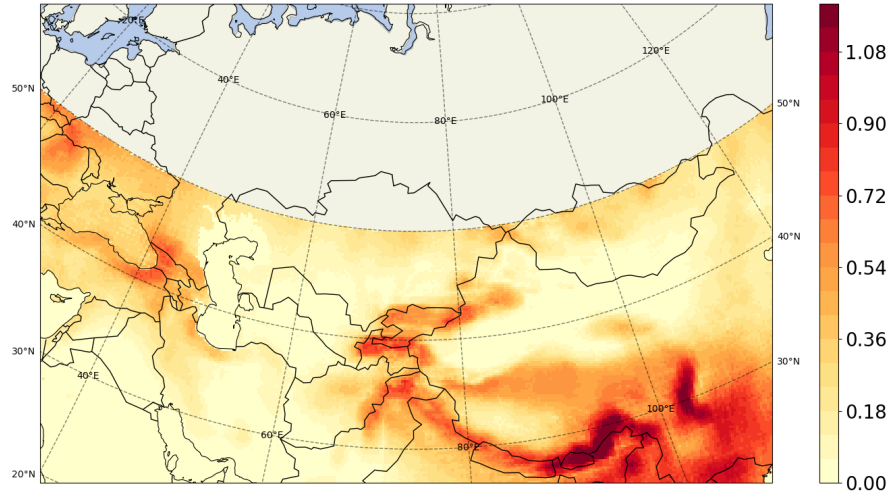
159 Relative probability differences and added values of ERAInterim reanalysis, CCLM, GERICS-REMO2015 and  
 160 RMIB-UGent-ALARO-0 driven by ERAInterim (hereafter, ERAInterim, CCLM-ERAInterim, GERICS-REMO2015-  
 161 ERAInterim and RMIB-UGent-ALARO-0-ERAInterim) are shown in Figure 2. The CHIRPS dataset is used as the  
 162 observational dataset  $O$  to calculate  $D$  according to equation 1. For example, D(ERAInterim) indicates the relative  
 163 probability difference between ERAInterim and CHIRPS. The measure of added value  $A$  is more pronounced over  
 164 areas with complex topography for all three RCMS (Figs.2.e-g). For example, the CCLM and GERICS-REMO2015  
 165 models show very negative values north of the Tibetan Plateau, whereas the RMIB-UGent-ALARO-0 model shows  
 166 positive values over corresponding areas. Over the Southeast of the domain, the CCLM and GERICS-REMO2015  
 167 show overall positive added values where the RMIB-UGent-ALARO-0 does not. Considering the whole domain,  
 168 all three models sensibly reduce the large and local-scale bias of ERAInterim against CHIRPS (Fig.2), especially  
 169 over complex topographies. Usually, the nested RCMs show similar values of  $D$  near their lateral boundaries, with



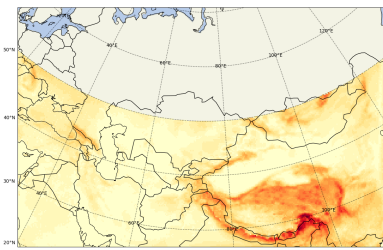
170 respect to ERAInterim (Fig. 2, panels *e,f,g*). However, even in those regions, e.g., in the South and Southeastern  
171 parts of the domain, positive values of  $A$  calculated with respect to CHIRPS are observed for the RCMs, but not for  
172 the GCMs. This might be because RCMs capture more detailed convective-forced precipitation than GCMs and we  
173 already removed 10 grid points of CCLM data at the lateral boundaries as the so-called buffer zone for our analysis  
174 where the model is forced to follow the forcing GCM. The precipitation pattern over the South and Southeast of the  
175 domain is mostly connected to the monsoon, and RCMs better fit the CHIRPS than ERAInterim (Fig. reffig:2.a–d).

#### 176 **Added value of CCLM driven by MPI-ESM1-2-HR1-2**

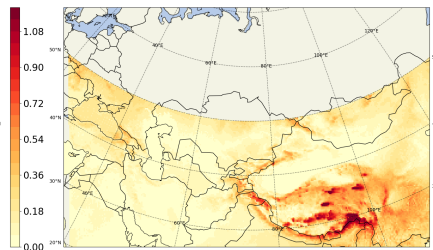
177 We showed that COSMO-CLM could reduce the large-scale bias of its driving reanalysis for daily precipitation,  
178 especially over areas with a complex topography and the Asian monsoon region. Here, we calculate the added value of  
179 the CCLM simulations driven by MPI-ESM1-2-HR for 1985-2014. It can be seen in figure 3.a that the driving GCM  
180 shows a more negligible bias than the ERAInterim over Tajikistan and Kyrgyzstan. According to Déqué et al. (2007),  
181 the GCM bias is one of the most important sources of uncertainty in the RCM’s regional climate projection, and  
182 the smaller  $D_{MPI-ESM1-2-HR2}$  compared to  $D_{ERAInterim}$  over Tajikistan and Kyrgyzstan might increase the skill  
183 of the final regional projections (under the assumption that the model bias remains conserved under other radiative  
184 forcings). The added value of CCLM driven by MPI-ESM1-2-HR shows small values over those areas compared to  
185 the simulation driven by ERAInterim. However, the bias over the South and the north band of the Tibetan Plateau  
186 and the East and West of the domain is reduced substantially. Our analysis of the two driving datasets (ERAInterim  
187 and MPI-ESM1-2-HR) tends to confirm the findings of the Sørland et al. (2018), at least of the total precipitation  
188 PDFs, that the biases of the GCM-RCM chain are not additive and not independent. For example, in almost all  
189 regions with high values of yearly precipitation, where GCM has a slight bias, the RCM does not present higher  
190 biases or vice versa.



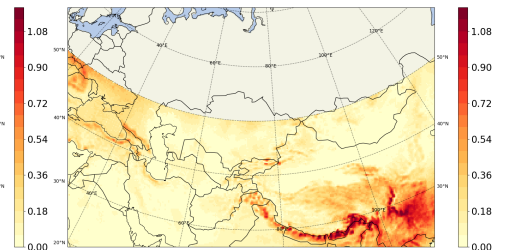
(a) D(ERAInterim)



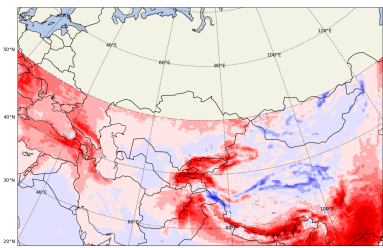
(b) D(CCLM-ERAInterim)



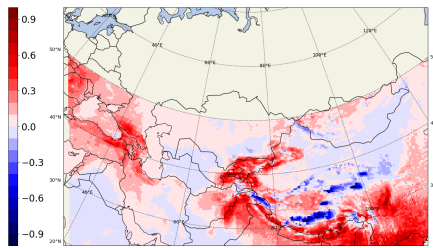
(c) D(GERICS-REMO2015-ERAInterim)



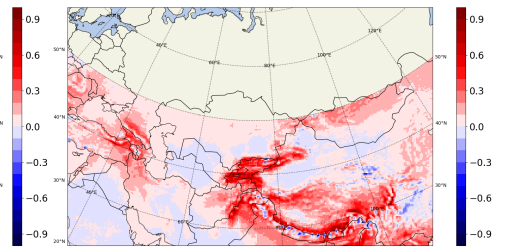
(d) D(RMIB-UGent-ALARO-0-ERAInterim)



(e)  $A_i$  CCLM



(f)  $A_i$  GERICS-REMO2015



(g)  $A_i$  RMIB-UGent-ALARO-0

Figure 2: Relative probability difference (D) for ERAInterim (a), CCLM-ERAInterim (b), GERICS-REMO2015-ERAInterim (c), ERAInterim-RMIB-UGent-ALARO-0 (d) and added value ( $A_i$ ) for the (e) CCLM-ERAInterim, (f) GERICS-REMO2015-ERAInterim and (g) ERAInterim-RMIB-UGent-ALARO-0 compared to the CHIRPS at 0.22° horizontal resolution.

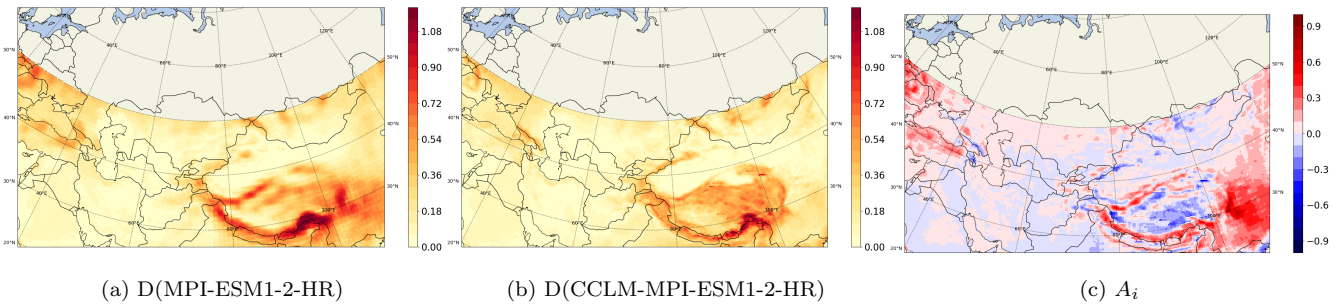


Figure 3: Relative probability difference (D) for MPI-ESM1-2-HR (a) and CCLM-MPI-ESM1-2-HR (b), and added value ( $A_i$ ) for the CCLM-MPI-ESM1-2-HR compared to the CHIRPS at  $0.22^\circ$  horizontal resolution.

191     **this should be f, i, l**

192     Done!

193     **P9: This doesn't make any sense when reading. So was the sentence before the highlight "areas**  
 194 **where ..."** about summer? And it looks like in winter there are actually big magenta blobs, especially  
 195 **for GERICS-REMO2015 (also for annual). You should also analyse this.**

196     We modified the whole two paragraphs in the new version of the manuscript according to your comment.

197     **P9: I'm guessing this comes from Fig.5? If yes it should be referenced. And again you use**  
 198 **"significantly" which is a very strong statement and requires statistical testing. Use a synonym.**  
 199 **Above when referencing you use Figs, should be Figs.4a-c for consistency... Also you're referencing**  
 200 **the GCM here while talking about the RCMs. I'd say remove the Fig reference and this should come**  
 201 **before the previous sentence talking about bias...This should come before and is not on a border, so**  
 202 **boundary effects don't seem to be a good explanation here What type of biases because one seems**  
 203 **to be positive (in green for tibetan) and the other negative (in magenta)**

204     We modified the whole paragraph accordingly. We made it clear that the RCM's performance near the boundaries  
 205 is influenced by the constraints of the driving GCM, reducing its ability to improve results in those regions.

206     **P11: No need to repeat this every time**

207     Done!

208     **Unless you actually compare the speed of dynamical DS versus ML that's an empty statement.**  
 209 **Same for computational cost.**

210     Done!

211 **P11: What's the time period of this analysis? It's also not mentioned in the Figure...**

212 Although the dataset covers a specific time period, the exact dates are not critical for this analysis because the  
213 data was randomly shuffled at the pair level before training, validation, and testing. This shuffling ensures that the  
214 model is exposed to diverse conditions without introducing temporal or spatial bias. Therefore, the time period is  
215 not relevant to evaluating the model's performance, as the focus is on ensuring generalization through randomized  
216 data exposure. For the distribution, 68,141 days (60%) were used for training, 22,714 days (20%) for validation, and  
217 22,714 days (20%) for independent testing, regardless of the specific dates.

218 While the dataset was shuffled for the original model training, validation, and testing to avoid temporal or spatial  
219 biases, ensuring the model's generalization performance, we took a different approach when evaluating the CNN with  
220 a new GCM later in the study. For this evaluation, the dataset was not shuffled to maintain the temporal structure  
221 and calculate the correlation, as the goal was to assess the model's performance with respect to the original time  
222 series of the new GCM. This allowed us to evaluate the model's ability to capture temporal patterns and validate  
223 the results accordingly.

224 We clarified that in the Figure.

225 **P11:it's relative to the RCM**

226 True, we modified the text.

227 **And what happens at the places where the constrained models reflect the GCM grid? Why is it**  
228 **negative there for the constrained models and not for the unconstrained?**

229 A very good point!

230 In Figure 4 of this answer we focus on the areas mentioned by you: green areas indicate improvement (reduction  
231 in error, i.e., positive added value) while magenta areas show worsening performance (negative added value). The  
232 blue lines highlight the regions where the GCM has maximum errors.

233 Interestingly, along these blue lines, the CNN shows no significant improvement compared to the GCM, implying  
234 that the CNN struggles to reduce the large errors where the GCM performs poorly. This behavior likely results  
235 from the constraints placed on the CNN—such as mass conservation or the preservation of physical quantities like  
236 precipitation—which limit its ability to freely adjust the outputs in these areas.

237 However, in the regions immediately surrounding these blue lines, we observe a mix of positive and negative added  
238 values, with green areas indicating that the CNN is successfully reducing errors, while the magenta areas suggest  
239 over- or under-compensation. This pattern may arise because, while the CNN is attempting to reduce the GCM's

240 large errors, it is constrained by the overall balance of physical variables (such as the mass of precipitation), causing  
241 it to redistribute errors into nearby regions.

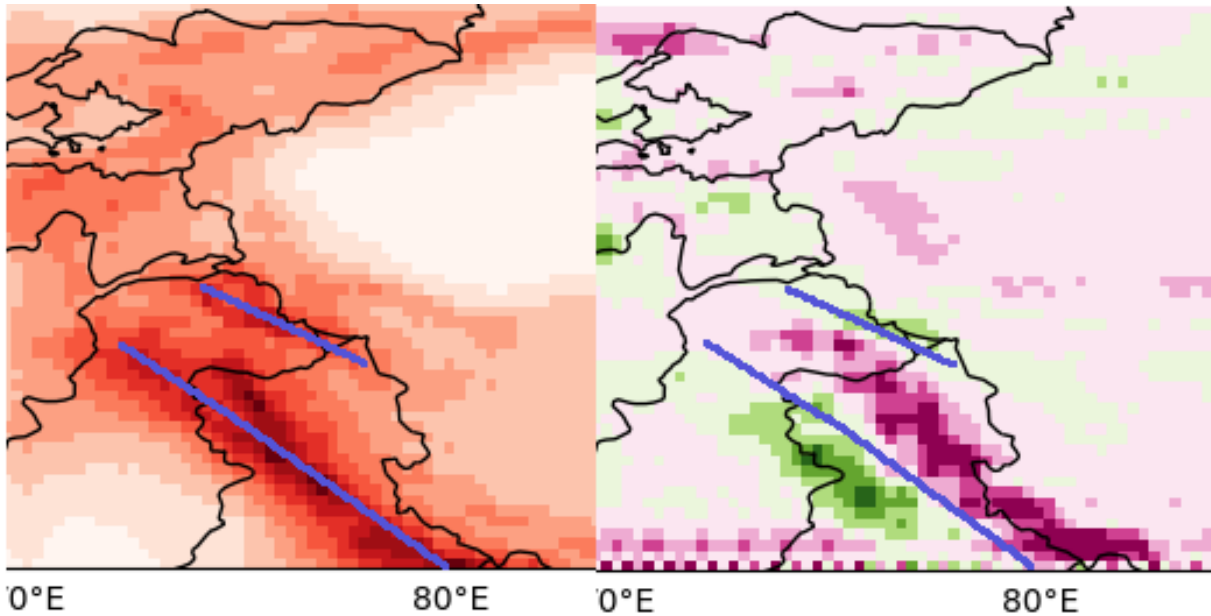


Figure 4: Left panel: GCM error distribution over the domain. Right panel: CNN error reduction compared to the GCM. Green areas represent regions where the CNN reduced errors (positive added value), while magenta areas show regions where the CNN increased errors (negative added value). The blue lines show places where the GCM shows maximum errors.

242 **Fig. for consistency... how much? This should be in the text... which is**  
243 Done!

244 **P12: You chose this GCM so it's not remarkable that they are very similar. You could have also**  
245 **chosen one that was different.**

246 Yes, we chose this GCM, but it was not cherry-picked. This model is one of the 10 selected models from the  
247 CMIP6 ensemble for the ISIMIP project. These 10 models were specifically chosen because they represent diverse  
248 physical processes and model codes across all components of the Earth system, ensuring a wide spread of the CMIP6  
249 models. The selection was designed to cover a broad range of potential outcomes, which is precisely why we used  
250 this model. For example Figure 5 of this answer shows how the precipitation annual trend for 1981-2014 are different  
251 in those models.

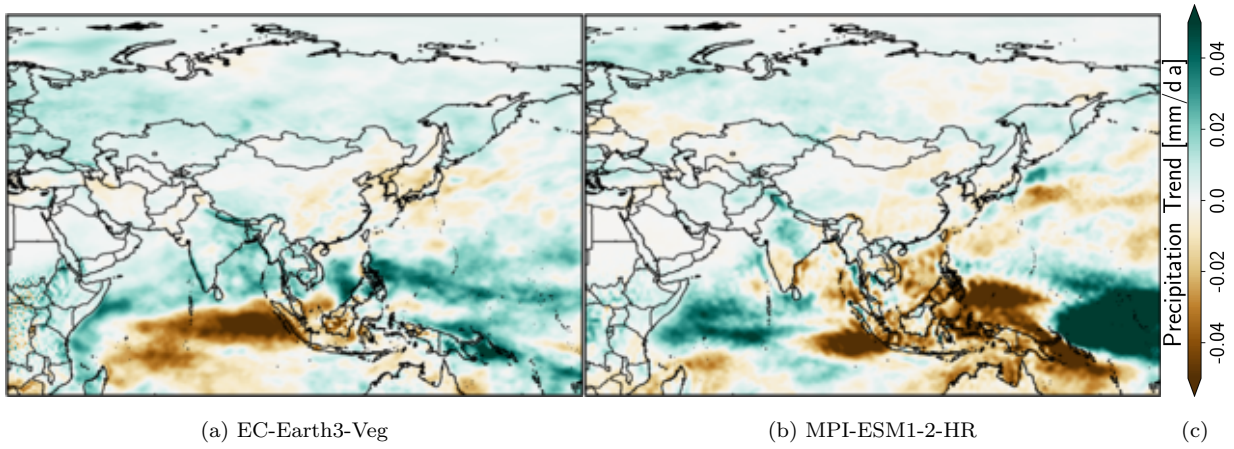


Figure 5: Asia Precipitation annual trend 1981-2014 single models of CMIP6

252 We also show the absolute values of pr then here in Figure 6.

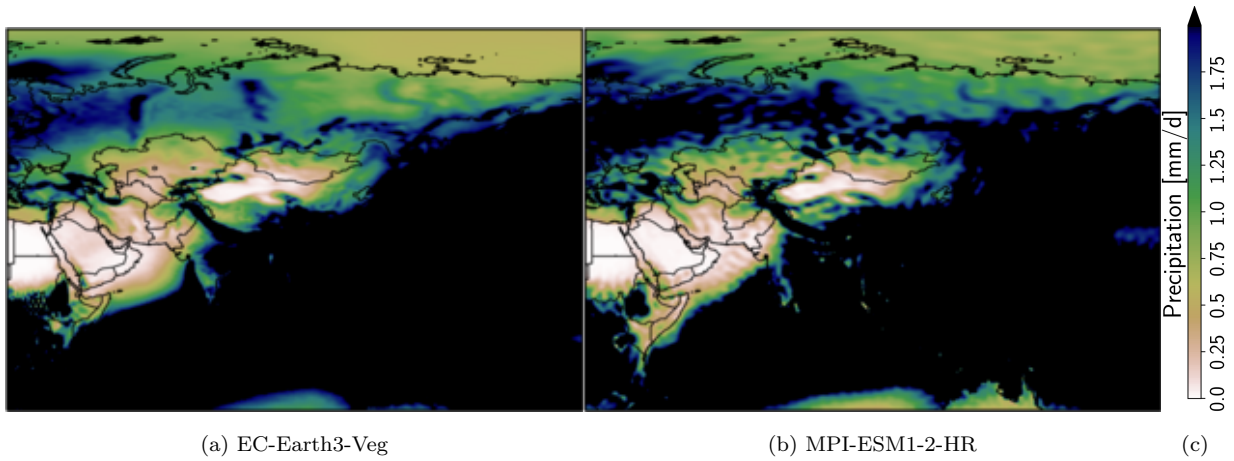


Figure 6: Asia Precipitation annual absolute 1981-2014 single models of CMIP6

253 But we remove the word "remarkable" in the new version of the manuscript. We agree that the generalization  
 254 ability depends on whether the new dataset is taken from the same distribution of values. They were driven by the  
 255 same forcing which are typical for any SSP scenario setups.

256 **P12: Yes but this was also not the case for the NoCL under the normal GCM so unsurprising.**

257 Done!

258 **P12: This should come before in the paragraph when the emulator is first mentioned.**

259 Done!

260 **P12: has improved compared to what? And give numbers..**

261 We added the reference and the numbers were already presented in the next sentence.

262 **P12: This should information come way before in the beginning of the paragraph as it's very**  
263 **important.. .**

264 We added that information to the beginning of the subsection and modified the text.

265 **P12: has a lower MAE than the GCM**

266 Done!

267 **P12: ressembles**

268 Done!

269 **P12: but how similar is this unseen forcing to the one the emulator has been trained on? This**  
270 **information is missing.**

271 Figure 7 shows the differences in total precipitation between the SSP scenarios for the MPI-ESM1-2-HR model.

272 As can be seen there are differences in the climatology of SSP370 and other SSPs.

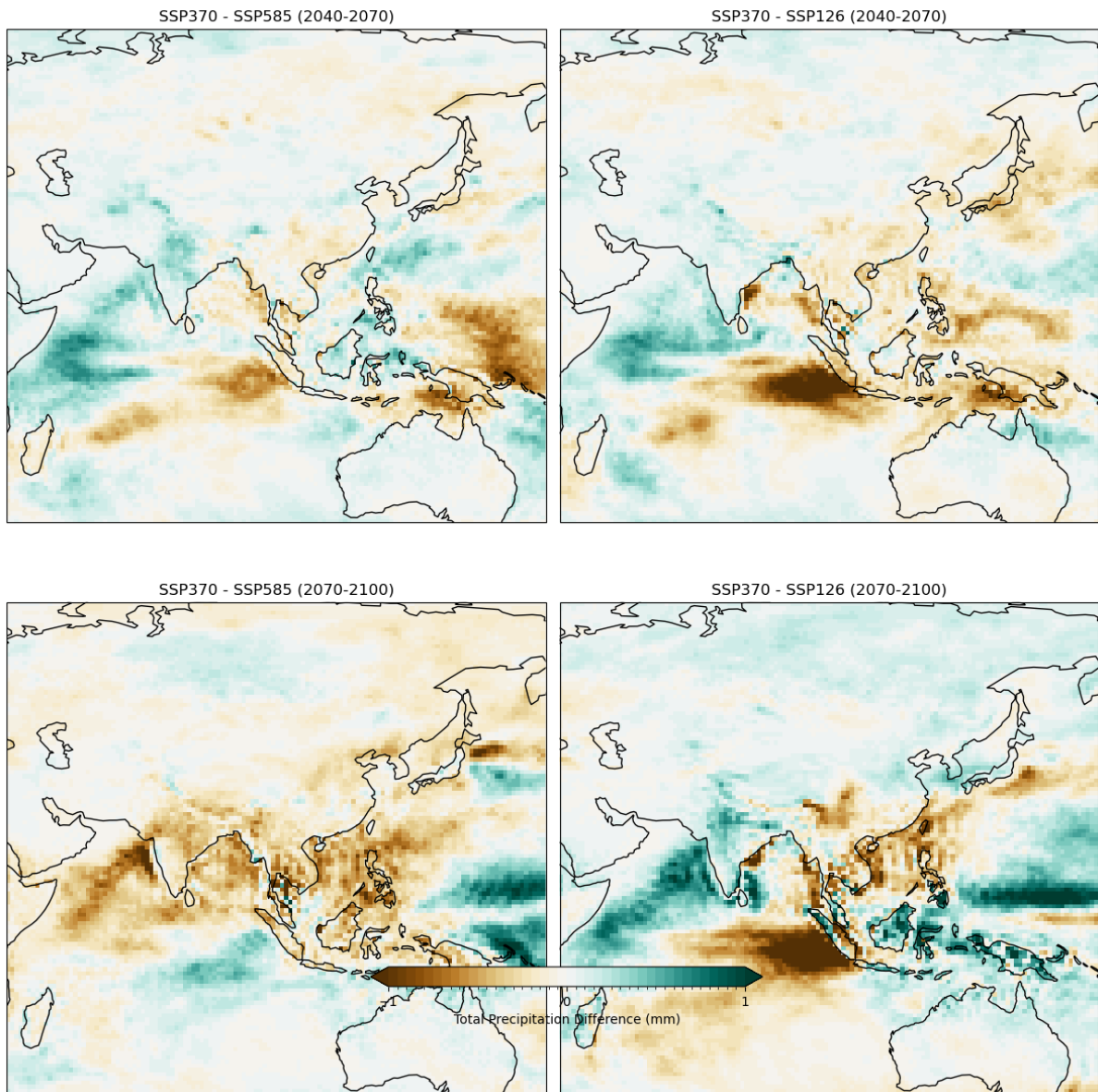


Figure 7: Differences in total precipitation between the SSP scenarios for the MPI-ESM1-2-HR model. The colorbar indicates the total precipitation difference in millimeters.

273 **P12:** No not really, it depends on how similar this new SSP is to the ones it was trained on,  
 274 this statement is coming on way too strong. Maybe rewrite into "shows promise for learning and  
 275 reproducing the [...].



276 We changed the text to what you have suggested.

277 **P13: of RCMs**

278 Done!

279 **P13: of RCMs**

280 Done!

281 **P13: Yes ok but you didn't do that so it feels weird to add this here. Except if you reformulate it**  
282 **into something that should be done in further work.**

283 We reformulated that to : "In future work, it would be valuable to follow the approach suggested by Volosciuk et al.  
284 (2017), where downscaling outputs are evaluated at coarser resolutions. This would allow for a deeper understanding  
285 of how downscaling methods introduce or fail to correct biases, which can vary significantly across spatial scales.  
286 By conducting evaluations on a coarser grid, we can better distinguish between the inherent biases of the model  
287 and those introduced by the downscaling process, providing important insights into the limitations and strengths of  
288 downscaling techniques in representing climatic variables across different scales."

289 **P13: No I really don't agree with this as Paleo climate experiments will be very different. Yes**  
290 **it learns a relationship but only for climate variables from the training set. So if you suddenly give**  
291 **it extreme values like for paleo climate it might not work and needs to be tested before making**  
292 **this statement. Guillaume Jouvét's CNN was a very different architecture and was trained on paleo**  
293 **climate directly not transferred like you propose. So this sentence feels wrong here.**

294 We removed those two sentences.

295 **P14: No again this statement is too strong. Maybe it "shows promise" but as I said before you**  
296 **didn't test the similarity of this unseen scenario to the training. Maybe it just fits nicely in there but**  
297 **it might not work for another more extreme scenario.**

298 We answered this and showed that the two models are not very similar. However, we modified the text in a way  
299 that it is not that strong.

300 **P21 : This figure ideally should also contain the time frame to know how many samples and from**  
301 **what time we're looking at.**

302 Done!

303 **P23 : MAE defined if you define added value.**

304 Done.

## 1.2 Other comments

Overall results: you use a lot of statements that need values from the figures to back them up; otherwise, they seem superficial. For example, if you say the "correlation of the CNN is higher than the RCM (here, include a value from the figures that supports this claim)." Also remove empty adjectives like "strikingly" or "remarkably" as they are very subjective and results should be based on numbers.

As already answered. We revised all those "superficial" adjectives and put numbers and analysis.

**3.1 Added value of CCLM driven by ERAInterim:** needs some shuffling of sentences, I would also present the MAE and bias results together instead of two separate paragraphs seeing how they are very similar. You can do this by using "this is also reflected in the bias blabla" somewhere.

Done in the new version of the manuscript.

**4 CCLM emulator using a CNN:** It needs values to support claims, and an analysis of the difference in negative AV for constrained CNN versus positive AV for unconstrained CNN would be nice (I'm thinking here of the magenta blob in the lower part of the image).

We discussed that before in this answer and added the analysis (Fig. 4 of this answer).

**4.1 Applying the CNN to a different GCM:** The new GCM you chose is very similar to the one the model is trained on. This should not be presented as a surprise but as a design choice that will impact the results, seeing how generalization is easier for the model then. I think that, in your case, generalization works pretty well, seeing how the map you produce is very similar to what you showed in Fig. However, the text should mention that for a more different GCM, it might work less well. The same applies to the new SSP; you should mention somewhere how similar this SSP is to what is in the training set. Otherwise, it's hard to evaluate the generalization ability. And if your model didn't work well on new GCM/SSP that is very different from training, it would not be a surprise or a downside of your model; that's just ML, but it needs to be transparent in the text, and that generalization ability depends on whether the new dataset is taken from the same distribution of values as the training and otherwise it needs special ML techniques like transfer learning.

In the revised manuscript, we have clarified that the new GCM we chose is not very similar to MPI-ESM. We included a comparison of the trends and climatology of these models and added a discussion on the ISIMIP model selection. Following your suggestion, we have removed any strong conclusions from the results and highlighted that

334 the performance of the CNN is indeed tied to the GCM it was trained on, as all models use the same SSP forcing.  
335 Additionally, we now emphasize the potential limitations of generalization when applying the model to more distinct  
336 GCMs or SSPs.

337 **Discussion: Following what I said above, I really disagree with the statement that your model could**  
338 **be applied to paleoclimate, as the distribution of paleo climate variables is probably very different**  
339 **from your training set. I would remove all this, also the mention of Jouvet’s work as it is very different**  
340 **from your case, and just mention that your CNN shows should work for RCMs from the same GCM**  
341 **as the training one, or for similar GCMs as you have shown in your results.**

342 Following your suggestion, we removed those parts from the new version of the manuscript.

343 **Generally, you show interesting results, but with your ML model, you should be very careful about**  
344 **the promises you make of its use and stay realistic, seeing how it is very limited to situations similar**  
345 **to its training (not because you did something wrong, but just because of ML in this setup).**

346 We completely agree with your point, and in the revised manuscript, we have carefully addressed this by ensuring  
347 that our promises regarding the use of the ML model remain realistic. We have taken care to present the model’s  
348 limitations and emphasize its applicability to situations similar to its training. .

---

349  
350 On behalf of all authors,

351 Bijan Fallah

## 352 References

353 Bengio, Y.: Practical recommendations for gradient-based training of deep architectures, in: Neural networks: Tricks  
354 of the trade: Second edition, pp. 437–478, Springer, 2012.

355 Ciarlo, J. M., Coppola, E., Fantini, A., Giorgi, F., Gao, X., Tong, Y., Glazer, R. H., Torres Alavez, J. A., Sines, T.,  
356 Pichelli, E., et al.: A new spatially distributed added value index for regional climate models: the EURO-CORDEX  
357 and the CORDEX-CORE highest resolution ensembles, *Climate Dynamics*, 57, 1403–1424, 2021.

358 Déqué, M., Rowell, D., Lüthi, D., Giorgi, F., Christensen, J., Rockel, B., Jacob, D., Kjellström, E., De Castro, M.,  
359 and van den Hurk, B.: An intercomparison of regional climate simulations for Europe: assessing uncertainties in  
360 model projections, *Climatic Change*, 81, 53–70, 2007.

- 361 Funk, C., Peterson, P., Landsfeld, M., Pedreros, D., Verdin, J., Shukla, S., Husak, G., Rowland, J., Harrison,  
362 L., Hoell, A., et al.: The climate hazards infrared precipitation with stations—a new environmental record for  
363 monitoring extremes, *Scientific data*, 2, 1–21, 2015.
- 364 Goodfellow, I.: *Deep learning*, 2016.
- 365 Harder, P., Hernandez-Garcia, A., Ramesh, V., Yang, Q., Sattegeri, P., Szwarcman, D., Watson, C., and Rolnick,  
366 D.: Hard-Constrained Deep Learning for Climate Downscaling, *Journal of Machine Learning Research*, 24, 1–40,  
367 2023.
- 368 Sørland, S. L., Schär, C., Lüthi, D., and Kjellström, E.: Bias patterns and climate change signals in GCM-RCM  
369 model chains, *Environmental Research Letters*, 13, 074017, 2018.
- 370 Volosciuk, C., Maraun, D., Vrac, M., and Widmann, M.: A combined statistical bias correction and stochastic  
371 downscaling method for precipitation, *Hydrology and Earth System Sciences*, 21, 1693–1719, 2017.