

# Manuscript NPG-2019-62

## Simulation-based comparison of multivariate ensemble post-processing methods

### Response to Reviewer Comments

*We thank the three reviewers for their positive assessment and their thoughtful comments which, we believe, will strengthen the manuscript. Below we addressed each comment in turn. Our replies are in italics. We also created a track-changes PDF for your convenience.*

*Kind regards,*

*Sebastian Lerch, Sándor Baran, Annette Möller, Jürgen Groß, Roman Schefzik, Stephan Hemri and Maximiliane Graeter*

#### Reviewer 1

This paper reports the results of a comprehensive simulation study comparing different methods for modeling spatial, temporal, and inter-variable correlations in statistically postprocessed ensemble forecasts. With a number of new multivariate methods having been developed in recent years, a study like this is of great interest as it allows readers to get an overview of the strengths and limitations of the different approaches.

General comments:

1. With the goal of this paper being a comparison between different methods for multivariate ensemble postprocessing, I feel that a more detailed discussion of the key features (e.g. optimal sampling of the predictive distribution vs. random sampling, assumption of stationarity of the copula structure vs. flow dependent copula structure, etc.) of the different methods should be given (possibly in the form of a table). This could serve as a motivation for the different simulation settings, which try to mimic situations where some of the assumptions are met while others are not.

*Thank you for this suggestion. We have added a table at the beginning of Section 2.2 where we now compare several key features of the different multivariate post-processing methods. We further refer to the related findings in Wilks (2015) (next comment) at the beginning of Section 2.2 and in the Discussion.*

2. Related to 1., I feel that the role of ensemble size (which has a big impact on the representation of the multivariate distribution) should be discussed a bit more. This seems relevant as for some methods it is easy to generate an ensemble of any size while for others it is not. I'm not suggesting that additional experiments should be performed, but a brief discussion of the findings in Wilks (2015) could be useful in a context where strengths and limitations of different multivariate postprocessing approaches are compared.

*We agree that the role of ensemble size is important and warranted a more extensive discussion. We have performed additional simulations for Settings 1 and 3 (Settings 1 and 2 in the revised paper) with ensemble sizes between 5 and 100. For Setting 1, the results are largely as it can*

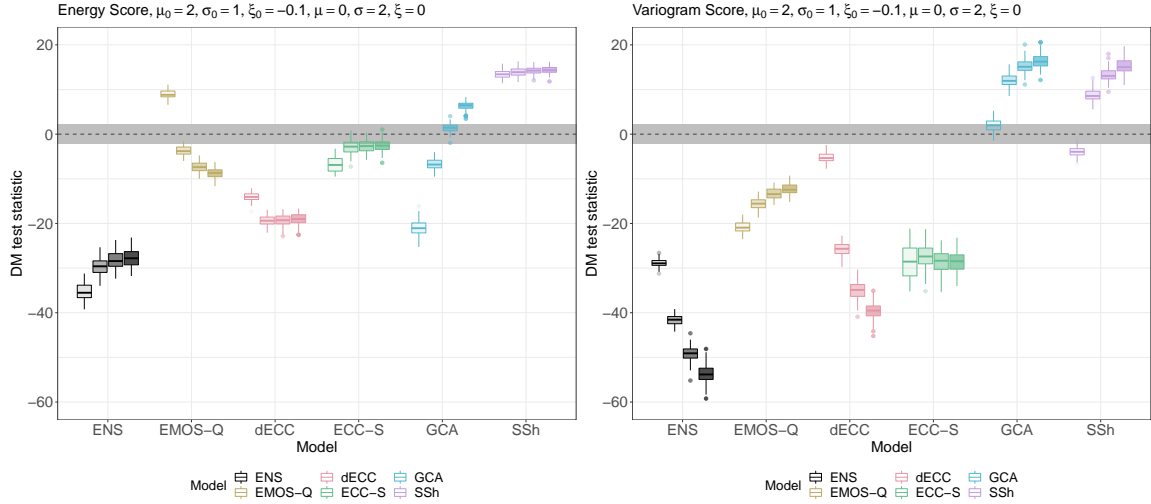


Figure R1: Effect of ensemble size in the GEV0 setting shown as grouped boxplots of ensemble sizes  $m = 5, 20, 50, 100$  (5 corresponds to lightest shade, 100 to darkest shade) for each model, for the additional Scenario with  $(\mu_0, \xi_0, \sigma_0) = (2.0, -0.1, 1.0)$  and  $(\mu, \xi, \sigma) = (0.0, 0.0, 2.0)$ , where  $\rho_0 = 0.75$  and  $\rho = 0.25$ .

be expected: The relative differences in terms of the ES between ECC-Q and ECC-S, and between ECC-Q and GCA become increasingly negligible with increasing ensemble size; likely due to increasingly smaller intervals from which random quantiles are drawn. Further, SSh shows improved predictive performance for larger numbers of ensemble members for  $\rho_0 < \rho$  in case of the ES, and for  $\rho_0 > \rho$  in case of the VS. This can likely be attributed to the corresponding increase in sample sizes when determining the dependence templates, similar to the effects on GCA. The relative performance of dECC is strongly effected by changes in  $m$  for large misspecifications in the correlation parameters. A positive effect of larger numbers of members relative to ECC-Q in terms of both scoring rules can be detected for  $\rho_0 > \rho$  when  $\sigma < 1$ , and for  $\rho_0 < \rho$  when  $\sigma > 1$ . In both cases, the corresponding effects are negative, when the misspecification in  $\sigma$  is reversed.

In Setting 2 (old Setting 3), in contrast to Setting 1, GCA seems to benefit most from an increasing number of members, while SSh benefits only slightly in terms of the ES, and stronger in terms of the VS. This is for example illustrated in Figure R1, which corresponds to the additional scenario presented in Figure 5 in the main paper. As in Setting 1, ECC-S becomes more similar in terms of the ES to the reference ECC-Q for increasing number of members. However, this effect, is not (or not that strongly) observable for the VS, where the number of members has nearly no effect on ECC-S.

We have added a paragraph to Sections 4.2.1 and 4.2.2 where the effect of ensemble size is discussed shortly. Figures with additional results have been added to the Supplemental Material in Sections 1.1 and 2.2 there. We are aware that displaying boxplots for several choices of  $m$  within a single figure is not optimal since the underlying random samples within each panel will necessarily differ across different values of  $m$ . Nonetheless, we believe that the differences due to random sampling are negligible due to the application of DM tests and the consideration of 100 repetitions of each simulation experiment. Therefore we chose the illustration added to the Supplemental Material in Sections 1.1 and 2.2 because the effects of changes in  $m$  are straightforward to compare.

3. Why has so much focus been given to simulation settings that are based on a time-independent model for the simulated forecasts and observations? I would argue that this (time-independence) is not a feature commonly encountered in applications, but with 3 out of 4 settings being

time-independent, multivariate methods that assume a stationary copula structure could be perceived as being more versatile than they really are. Agreed that setting 1 is a natural starting point for such a comparison and that setting 3 is interesting because of the entirely different nature of the marginal distributions (skewness, possibility of heavy tails, mixed discrete-continuous distributions), but what do we learn from setting 2 that we cannot learn from 1 and 3? The main difference to 1 seems to be the poorer performance of GCA, but an explanation for this is not given, and so the insights gained from this setting are limited.

*We agree that the results of Setting 2 were overall very similar to those of Setting 1. Therefore, we have removed Setting 2 from the paper and added a slightly adjusted version to the Supplemental Material. The former Setting 2 is there denoted by Setting S1, and setup and results are now discussed in Section 5 of the Supplemental Material. Accordingly, the paper has been adjusted as follows:*

- *Setting 4 is now Setting 3, and Setting 3 is now Setting 2*
- *Sections 3.2 and 4.2.2 have been removed from the paper and added to the Supplemental Material*
- *Figure 1 has been adjusted by removing results for Setting 2*
- *several references to Setting 2 in the text have been removed throughout*
- *A short paragraph referring to the additional setting (now in the Supplemental Material) has been added to the end of Section 3.1.*

Specific comments:

131-132: Please check if that statement is correct. In my recollection the selection of past observations in Clark et al. (2004) was not random, but was based on the valid date of the forecast

*Thank you for pointing this out, you are of course correct. The description of SSh has been corrected and information about the random selection of training cases has been added to Section 3.1 in step (S3).*

293-297: I find setting 4 the most interesting, but I find the particular definitions of the model parameters unnecessarily complicated. Specifically, I don't get a good intuition of what kind of time-varying correlations this model implies. Couldn't one simply define

$$\Sigma_{i,j} = \sigma \rho^{|i-j|}$$

as in the other settings, but now make  $\rho$  time-varying, e.g. via

$$\rho(t) = \rho_0(1 - a/2) + \rho_0(a/2) \sin\left(\frac{2\pi t}{n}\right).$$

This model would be autoregressive with lag-1 correlations oscillating between  $\rho_0$  and  $\rho_0(1 - a)$ , and thus have a more intuitive interpretation.

*Thank you for the suggestion of this alternative setting. We agree that this definition has a more intuitive interpretation. However, we believe that the way in which  $\rho$  and  $\rho_0$  are defined does not directly correspond to the situation we aimed to cover in Setting 4: Our goal was to mimic a situation in which the covariance structure of both observations and ensemble varies over time, with possible misspecifications of this structure in the ensemble. In the setup suggested above, the covariance structure of the observations does not change over time. Therefore, non-time-varying methods such as GCA and SSh which assume stationarity of the covariance structure are*

at an advantage in that they do not suffer from the (time-varying) misspecifications in the raw ensemble. An exemplary illustration is given in Figure R2. Note that the rows show results for different values of  $a$ . SSh here never performs worse than ECC-Q and all methods significantly outperform ECC-Q in terms of the VS for almost all parameter values.

In order to provide a time-varying simulation with a possibly more intuitive interpretation, we have added a variant of the setting suggested above to the paper. In this new Setting 3B, we set  $\Sigma_{i,j}^0(t) = \rho_0^{|i-j|}(t)$ , for  $i, j = 1, \dots, d$ , where the correlation parameter  $\rho_0(t)$  varies over iterations according to

$$\rho_0(t) = \rho_0 \cdot \left(1 - \frac{a}{2}\right) + \rho_0 \cdot \left(\frac{a}{2}\right) \sin\left(\frac{2\pi t}{n}\right)$$

for  $a \in (0, 1)$  for the observations, and similarly,  $\Sigma_{i,j}(t) = \sigma \rho_{|i-j|}(t)$ , for  $i, j = 1, \dots, d$ , where

$$\rho(t) = \rho \cdot \left(1 - \frac{a}{2}\right) + \rho \cdot \left(\frac{a}{2}\right) \sin\left(\frac{2\pi t}{n}\right)$$

for the ensemble predictions. Correlations in both cases are autoregressive with lag-1 and oscillated between  $\rho_0$  and  $\rho_0(1 - a)$ , and  $\rho$  and  $\rho(1 - a)$ , respectively.

The description of Setting 3B was added to Section 3.3, results are discussed in Section 4.2.3, and a corresponding figure has been added to that section. Results for additional simulation parameters for Setting 3B are provided in the Supplemental Material.

260-264: I find this notation a bit confusing since previously the subscript/superscript 'O' was used for observations and here the subscript '0' (which is hard to discern from 'O' in the NPG font) is used to denote the fraction of zero values. The notation is also inconsistent in that in setting 3 'x' and 'y' are used to denote forecasts and observations, in contrast to the subscript/superscript 'O' for observations in the other settings.

The subscript '0' (to denote zero values) has been changed to subscript 'z'. The subscripts 'x' and 'y' have been changed to match the subscript notation in the other sections.

323 '... are identical to those of ECC-Q ...': Is this really true for ECC-S? The way it is described here, ECC-S seems to imply some level of randomization (albeit less than ECC-R), so the sampling is not the same as for ECC-Q.

Thank you for spotting this. The univariate distributions of ECC-S will indeed not be identical to those of ECC-Q. Comparing univariate results, we however found that the differences in terms of predictive performance are only very minor and were not noticeable in plots. We have modified the corresponding paragraph which now reads "Note that for ECC-S and SSh differences in the univariate forecast distributions compared to those of ECC-Q may arise from randomly sampling the quantile levels in ECC-S and due to random fluctuations due to the 10 random repetitions that were performed to account for simulation uncertainty of those methods. However, we found the effects on the univariate results to be negligible and omit ECC-S, dECC and SSh from Figure 1."

Fig. 2: I don't think that this figure is really necessary. Why is the (univariate) performance of ECC-Q compared to the raw ensemble relevant to the comparison of multi-variate postprocessing approaches?

Figure 2 and the corresponding text paragraph have been removed.

354 '... SSh never performs substantially worse ...': Why would we expect otherwise? The only drawback of SSh in the present context is the underlying assumption of time-invariance of the correlation structure, which is not a drawback in a time-invariant simulation setting. If not discussed before, this paragraph could be an opportunity to discuss this issue of time-invariant vs time-varying correlations.

We have added a short discussion to the corresponding paragraph.

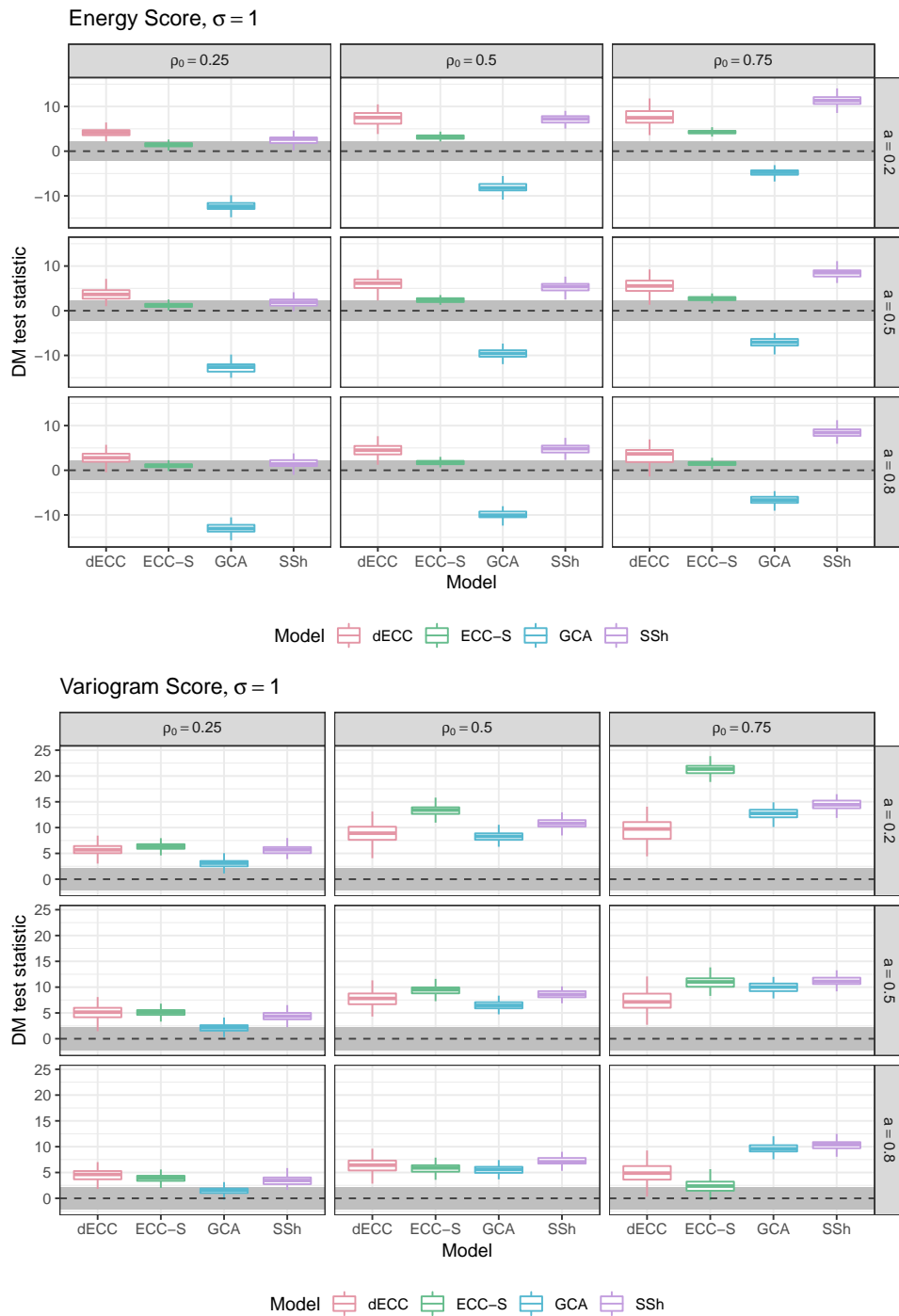


Figure R2: Illustration similar to the figure for Setting 4 in the original main paper, but based on the adapted simulation setting suggested by Reviewer 1.

360: I wonder if ECC-S gives the better results for the wrong reason here. Maybe I am misunderstanding its key idea, but to me this is essentially a compromise between ECC-Q and ECC-R. Is it possible that the small amount of randomization in ECC-S weakens the correlations, which is beneficial for  $\rho > \rho_0$  and detrimental for  $\rho < \rho_0$ ? An argument against this hypothesis is that the performance of ECC-S is not significantly different from ECC-Q when  $\rho = \rho_0$ . I just cannot think of any good reason why ECC-S would be better than ECC-Q. If the authors have any explanation for these results I would encourage them to include those in the discussion of the results.

*Unfortunately, we are unable to provide a comprehensive explanation for these observations. As will be argued in the response to Reviewer 2 below, drawing conclusions on specific aspects of the multivariate methods or on the effect of single simulation parameters is difficult and is further impeded by a lack of well-understood multivariate evaluation methods. One potential advantage of the sampling scheme in ECC-S in comparison to the use of quantiles at fixed levels can become apparent when considering corresponding EMOS-S and EMOS-Q variants. When evaluated in terms of the VS, the random sampling in ECC-S may alleviate issues arising from over-estimated correlations due to the use of the same quantile levels in EMOS-Q. However, it appears to not be straightforward how these observations will be effected by applying ECC, and what the effects of the 10 repetitions of each individual simulation experiments that were performed to account for simulation uncertainty, would be.*

Fig. 6: The caption should state that this is scenario B from setting 3

*Fixed.*

458 'changes over iterations': I would change this terminology and speak of 'time' instead of 'iterations', and be more specific about what changes/varies over time. Like-wise in 464-465 I would clarify what you mean by 'structural change'

*Changed as suggested.*

Language and typos:

43: ... studies allow one to specifically tailor ...

*Fixed.*

135: sufficiently many

*Fixed.*

137: 'not directly straight forward' sounds weird, I'd just say 'not straightforward'

*Changed as suggested.*

249: ... allows one to generate ...

*Fixed.*

375: 'can potentially be explained' sounds weird, maybe better 'may be explained'

*Changed as suggested.*

388: I think you want to say 'In terms of ...'

*Fixed.*

443: Change to 'the other way round' and 'In accordance with'

*Changed as suggested.*

445: Change to 'in contrast to'

*Changed as suggested.*

## Reviewer 2

In this paper, a comprehensive simulation study is implemented, comparing multiple multivariate statistical post-processing methods which all combine standard univariate post-processing with one of four techniques to reintroduce spatial, temporal or inter-variable dependencies. It is well-written, an important contribution given the variety of techniques available and potentially very useful to identify optimal operational post-processing strategies for varying types of data. Just a few clarifications and changes are needed.

General comments:

I would have liked to have more focus (or at least comments) on the effect of ensemble size and dimension. The here chosen ensemble of 50 members is at the upper limit of what is now operationally produced, with the majority far below this number. Of course it is important for the study to have a sufficient number of data points so as to produce significant results, but it would also be interesting to look at settings with a smaller number of ensemble members. Also, the number of dimensions ranges from 4 to 5, which would correspond to looking at consistency between a few weather variables, but would usually be too low for a setting where preserving spatial or temporal features are important. I wonder if the findings would be different for smaller ensembles or higher dimensions.

*Thank you for this suggestion. Following a similar comment by Reviewer 1, we performed additional simulations to assess the effects of ensemble size. Please see the response to Reviewer 1 for a discussion of the role of the number of ensemble members.*

*To study the effects of the number of dimensions, we performed additional simulations for Setting 1 with dimensions between 2 and 50. Overall, the relative results are often not effected too much by changes in the number of dimensions, in particular in terms of the energy score. Somewhat more substantial differences can be observed in terms of the variogram score. In a nutshell, GCA performs worse for higher dimensions, whereas ECC-S improves for larger values of  $d$ . The relative differences to SSh (in favor of SSh) become more substantial with increasing  $d$ , whereas the changes in relative performance of dECC show a strong dependency on the combined misspecification of variance and correlation of the raw ensemble.*

*We have added a paragraph to Section 4.2.1. Additional figures with results for  $d = 2, 3, 10, 20, 30$  and 50 have been added to the Supplemental Material (Section 1.2 there). See the response to the corresponding comment by Reviewer 1 for a discussion of the display as multiple boxplots within one figure. Performing additional simulations to investigate different choices of  $d$  for Setting 3 (which is Setting 2 in the revised manuscript) was impossible due to constraints regarding the computational requirements, see also our comment below.*

I find it very interesting that the performance of certain methods is sometimes very different when  $\rho > \rho_0$  than in the opposite case. Do you have any explanation for this?

*We believe that despite the (relative) simplicity of the chosen simulation setup, interpreting differences in multivariate performance is challenging due to inter-connected contributions of misspecifications in mean, variance and correlation structure and their effects on the (still) not well understood multivariate proper scoring rules. We are thus unable to provide a comprehensive explanation for the particular role of the correlation parameters  $\rho$  and  $\rho_0$ . Some potential explanations of differences in forecast performance may be given by observing that under certain circumstances, modifying a raw ensemble prediction by artificially weakening correlations, for example by randomly permuting ensemble member's forecast vectors, may unexpectedly improve predictive performance. This was for example observed in Schefzik (2017), and is likely an issue of underdispersed univariate forecast distributions that improve by the changes on the univariate forecast distributions imposed by the modifications of the correlation structure. However, to*

*assess the effects of different sources of misspecifications in further detail, additional multivariate verification techniques such as multivariate rank histograms should accompany the analysis. While such additional studies are beyond the scope of the paper, they represent an interesting starting point for future research.*

#### Specific comments

1. Line 71 and Line 153: The correlation matrix here is not necessarily the identity matrix, so I don't think it is a standard normal distribution.

*Fixed.*

2. Line 74: I would mention here that  $m$  is the number of ensemble members

*Added.*

3. Section 2.2: There is a mixture of  $x$  and  $X$  used to define samples and ensemble forecasts, but I was confused why this distinction is made within the notation, it seems inconsistent.

*Throughout the description of the methods, we used  $x$  to denote univariate quantities in the individual dimensions.  $X$  in bold print is used to represent vector-valued quantities, and  $X$  in normal print is used to for components thereof. A sentence has been added at the beginning of Section 2.2 to clarify this.*

4. Lines 184-186: For the other settings it is mentioned to which weather variables these settings could apply. It would be nice to add something like this to Setting 1, as well.

*A corresponding sentence has been added to the beginning of Section 3.1.*

5. Line 242: The notation in this setting is different from the others and this is confusing. Here, the forecasts and observations are marked with  $x$  and  $y$ , whereas the other settings use  $o/0$  to mark the observations.

*The notation has been changed accordingly, see also our answer to Reviewer 1.*

6. Line 276: Is there a specific reason, why  $d$  is 4 in this setting and 5 in the others?

*The limitation to  $d = 4$  was mainly due to computational requirements and numerical stability issues caused by the NORTARA package for R used to generate the multivariate ensemble predictions and observations. A sentence has been added in the paper.*

7. Lines 292-293: Some of the matrices are in bold face, some are not.

*Fixed. Additional minor corrections to the description of Setting 4 (now Setting 3A) are detailed under "Further changes".*

8. Figure 1: I am a bit surprised to see that GCA is performing that much worse as compared to the other post-processing methods (in a univariate sense). There are even cases where the performance is equal or possibly worse than for the raw ensemble. Do you have an idea why that could be?

*In particular when compared to ECC-Q and related methods, we believe that this is mostly an issue of sampling. As discussed in Section 4.1, the univariate quantile forecasts of ECC-Q are (close to) optimal in terms of the CRPS, but the random samples from the predictive distributions obtained in GCA are not (see next comment). We have modified the corresponding sentence to make this more clear. Cases where performance is worse than the raw ensemble are very rare and likely arise when simulation parameters of the ensemble forecasts are close to these of the observations.*



9. Lines 328-330: Can you explain a bit further what you mean by “optimal in the terms of the CRPS”?

*When the CRPS is used for evaluating a forecast distribution represented by a finite simulated sample, the sample is implicitly interpreted as a set of quantiles from the underlying forecast distribution at levels  $\frac{i-0.5}{m}$ ,  $i = 1, \dots, m$ . Viewed differently, this implies that utilizing quantiles (at these levels) when generating a univariate sample from a forecast distribution will result in a lower expected CRPS than drawing samples at random. Even though the levels at which quantiles are obtained are not identical to the optimal quantile levels, the differences will be small for  $m = 50$ . Details on the mathematical background can be found in the referenced paper by Bröcker (2012). We hope that this becomes more clear with the modifications made to this part of the paper mentioned in the response to the preceding comment.*

10. Lines 334-335: Naturally, scenario D has the smallest improvement compared to the others. Does that also mean that the scenarios are on the same absolute skill level after post-processing?

*We did not further investigate univariate performance since we removed this paragraph following a suggestion of Reviewer 1.*

11. Footnote 4: In my opinion it would be clearer if you refer to ECC-Q as EMOS-Q in this section as well.

*We have modified the corresponding paragraph. Together with the new table comparing key features of multivariate post-processing methods (see comment by Reviewer 1), we hope that this sufficiently clarifies terminology.*

12. Lines 415-430: Can you refer to the figures in the appendix that show these results by number?

*References to the corresponding sections of the Supplemental Material have been added.*

13. Line 446: “the VS might be better able to account..” This is confirming a known result, therefore “might” is a bit unsuitable.

*Changed as suggested.*

#### Technical corrections

1. Line 327: I would move the sentence beginning ”Note that” to footnote 4, as it directly relates to the changes in the marginal distributions mentioned there.

*The corresponding paragraph has been modified following comments from Reviewer 1, and the footnote has been removed.*

2. Line 356: I find this sentence a bit confusing. Should there be a comma before “the less information”?

*We have modified the sentence and added a comma as suggested.*

3. Line 386: Missing comma before “where”

*Fixed.*

4. Lines 415 and 441: I would add “parameter” after “observation location”.

*Changed as suggested.*

### Reviewer 3 (Zied Ben Bouallègue)

This paper presents results of an intercomparison study. Different multivariate ensemble post-processing methods are compared using toy-model simulations. The focus is on empirical copula methods that are generally applied as a second post-processing step, after univariate post-processing step, in order to provide coherent multivariate structure to ensemble calibrated forecasts.

As the reviewer is the developer of one of the methods compared here, he prefers to make himself known. There exists no conflict of interest (*stricto sensu*) but potential cognitive biases from the reviewer side when scrutinizing the results. The paper is clear, well structured, and well written. However, the choice regarding the selection of illustrations is not sufficiently motivated to my opinion. This choice is important because it drives the discussion and the main conclusion of the study. A 3-point argument is developed below to explain this criticism.

You claim in the conclusion (L462/463) that the 4 simulation settings aim to mimic different situations and challenges occurring in practical situations. However, the link with practical situations is sometimes weak or missing. In particular, it would be interesting to link misspecification definitions with practical examples. What  $\sigma > 1$  and  $\sigma < 1$  mean, and similarly what does  $\rho > \rho_0$  and  $\rho < \rho_0$  mean and “look like” in practice? In which situations should one expect to encounter these types of misspecification? Which type of misspecification situations are the most common in practice? Are combinations of misspecified elements more common than others ( $\sigma < 0$  and  $\rho < \rho_0$  for example)?

Once the link to the applications is clarified, illustrations could be chosen consistently. Result material is abundant, so one selection criterion could be to focus on the main misspecification encountered in practical situations. For example, you use  $\sigma = 1$  to illustrate results in Setting 2. Does that often occur in practice? One could rather illustrate Setting 2 with  $\sigma \neq 1$ . Similarly, for Setting 3, scenario B is used for illustration purposes. It corresponds to the case where the ensemble forecasts have a heavier right tail and slightly higher point mass at zero than the observations. Does that often occur in practice? No justification for the use of this scenario as reference is provided. Scenario A (the ensemble comes from a distribution with smaller spread) or Scenario C (the observation distribution has a much heavier right tail) seem to be more likely to be faced in practice.

*Thank you for this helpful comment which sparked ample discussion between the authors about how to best diagnose and transfer misspecifications of real-world ensemble prediction system to the simulation settings.*

*From our experience with most standard surface weather variables, one would expect that the univariate ensemble predictions are typically underdispersive ( $\sigma < 1$ ), and exhibit a bias ( $\epsilon \neq 0$ ). Often, biases and dispersion errors will of course vary over time, and may be vastly different for different variables or geographical locations. The choice of the correlation parameters  $\rho, \rho_0$  naturally hinges on the chosen correlation function for the simulation setting. In practice, this would likely correspond to an over-simplification of true correlations. In an attempt to diagnose realistic correlation parameters in the context of Setting 1, we have estimated correlation parameters for 2-day ahead ECMWF ensemble predictions of 00UTC temperatures at observation stations in Germany based on the 10-year dataset used in Rasp and Lerch (2018), as follows: For a randomly selected station, we choose the 13th, 26th, 38th, 50th closest station. If there are no substantial differences in altitude, we determine the empirical correlation among the observations at those 5 stations, as well as the correlation among each ensemble member’s forecasts at those stations. Next, the differences between the estimated correlation in the ensemble members and the observation are computed for all 50 members. In addition, we determine the parameters  $\rho, \rho_0$  from the exponential correlation function model assumed in the paper by numerical optimization. The procedure described above is repeated 100 times. Figure R3 shows differences in*

the empirical correlation, with histograms summarizing results over all 50 members and 100 repetitions. The differences in correlation are almost identical for all close stations, suggesting that the assumption of a fixed parameter  $\rho$  and  $\rho_0$  seems reasonable. Further, a similar conclusion can be obtained from Figure R4 which shows differences in the estimated correlation parameter of the exponential correlation model. Both figures suggest that correlations in the observations are over-estimated by the ensemble. Realistic settings thus probably relate best to simulation settings where  $\rho > \rho_0$ , but the values chosen in the paper (with differences of at least 0.15) appear to lead to possibly be too large differences (at least when compared to 2-day ahead temperature predictions over Germany).

In deciding on parameters for the simulation settings, we have mainly sought to cover a complete range of possible misspecifications rather than mirroring the situation in practice, in order to provide a more complete view of the performance of the multivariate post-processing methods. We have extended the discussion on the realism (or lack thereof) of the simulation settings in the Discussion, and have incorporated some of the arguments made above. Further, realistic values of the simulation parameters of Setting 1 that can be expected in practical applications are now discussed towards the end of Section 3.1. Setting 2 has been removed following the suggestion of Reviewer 1. Concerning the interpretation of the chosen scenarios in (former) Setting 3, for example Scenario B considers a situation where the forecast distribution estimates the amount of zero precipitation correctly, but otherwise the probability for obtaining a value smaller or equal to a fixed precipitation amount  $x$  computed from the forecast distribution is always smaller than the corresponding probability computed from the distribution of the observation, see the right panel in Figure R5. In combination, these two features may not occur very often in practice, since one would expect that underforecasting smaller precipitation amounts should come along with an underestimation of zero precipitation. See the left panel in Figure R5, showing fitted CDFs to a sample of observed precipitation values and corresponding forecasts of an individual ensemble member at a specific station in Germany based on real precipitation and ECMWF ensemble forecast data. Since a considerable number of scenarios with respect to the actual values of  $\mu$ ,  $\sigma$  and  $\xi$  is conceivable in practice, depending on the climatological circumstances, further investigations based on real data are required to provide additional insights. Furthermore, the interplay between the 3 GEV0 parameters is specifically complex in the sense that all 3 parameters have a joint influence on the location and dispersion properties, so that simple misspecifications in mean and variance might correspond to various (different) sets of parameter combinations. However, we agree that a more detailed investigation of theoretical and practical properties of the GEV0 distribution is a highly interesting starting point for future research.

Your conclusion points to the robustness of the SSh method and so you encourage post-processing practitioners to consider this method as their first choice. Based on the results presented in the manuscript and the ones in the supplemental material document, one could draw the opposite recommendation. First choice methods are available for different types of misspecification and so one could encourage to apply SSh only when the misspecification is unknown or difficult to identify. Specific recommendations would read:

1. If the ensemble is underdispersive and the ensemble correlation is too weak, d-ECC is the best option, both in terms of ES and VS.
2. If the ensemble is underdispersive and the ensemble correlation is too strong, ECC-S is the best option in terms of VS and one of the best options in terms of ES.
3. If the ensemble is overdispersive and the ensemble correlation is too strong, d-ECC is the best option, both in terms of ES and VS.
4. If the ensemble is overdispersive and the ensemble correlation is too weak, GCA is the best option in terms of VS, but is less performant in terms of ES. GCA proves to work well even when the overdispersiveness error characteristic is relaxed.

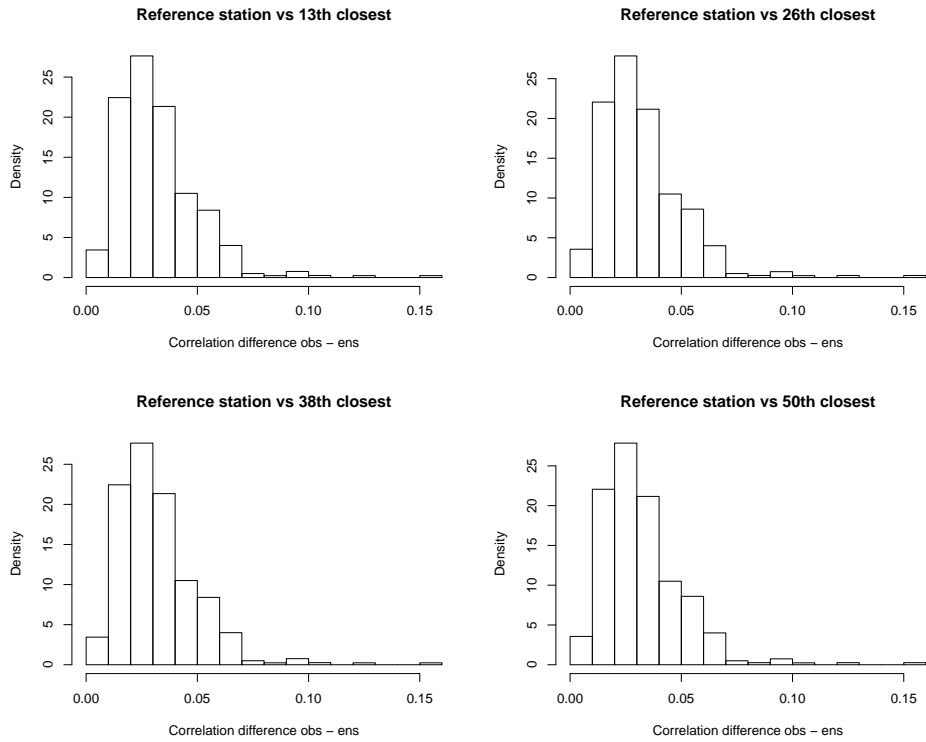


Figure R3: Differences in empirical correlation of observations and ensemble members at a set of 5 close stations based on the dataset of Rasp and Lerch (2018).

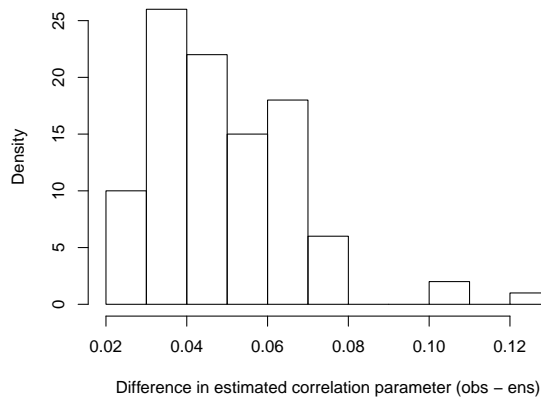


Figure R4: Differences in estimated correlation of observations and ensemble members according to the exponential correlation function model assumed in the simulation settings at a set of 5 close stations based on the dataset of Rasp and Lerch (2018).

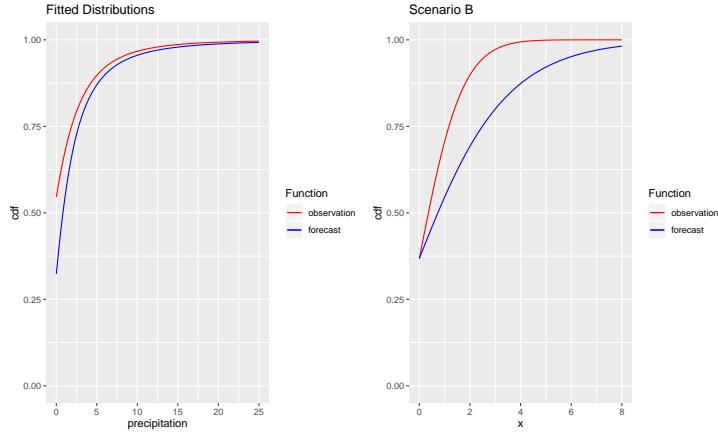


Figure R5: Cumulative distribution functions of  $GEV_0$ . The fit parameters in the left panel are  $\mu_0 = -1.0698$ ,  $\xi_0 = 0.2700$ ,  $\sigma_0 = 1.9906$  for the observation (red line) and  $\mu = 0.1887$ ,  $\xi = 0.4063$ ,  $\sigma = 1.5890$  for the forecast (blue line). The fit parameters in the right panel are  $\mu_0 = 0$ ,  $\xi_0 = -0.1$ ,  $\sigma_0 = 1$  for the observation (red line) and  $\mu = 0$ ,  $\xi = 0$ ,  $\sigma = 2$  for the forecast (blue line).

5. Otherwise consider using SSh, in particular when the correlation structures in the forecasts and observations are very dissimilar.

This is valid for all types of distributions (so all types of weather variables). Are these conclusions still valid in case of time varying misspecification? Verification results are missing to conclude here. Therefore, the authors are encouraged to investigate various sigmas in Setting 4 in order to collect evidences, and could draw conclusions accordingly

*Thank you for these suggestions. We agree that the conclusions may benefit from a more detailed discussion of situations where the different methods show advantages or disadvantages across settings.*

*With the changes and additional results in the paper and Supplemental Material following the comments by Reviewer 1 (Setting 2 was removed; Setting 4 (now Setting 3A) is accompanied by another variant, Setting 3B; and additional simulations with varying numbers of ensemble members and dimensions), we believe that it is difficult to arrive at general conclusions of this form that would be valid for all types of distributions and settings. In particular, the effects of misspecifications of the individual simulation parameters may be very different in, for example, the multivariate Gaussian distribution in Setting 1 and the censored GEV variant in Setting 3 (now Setting 2). In particular for the new time-varying setting, some of the results regarding the relative performance of dECC and ECC-S, respectively, differ somewhat from the recommendations summarized above. With Setting 1 in mind, realistic parameter choices to mimic properties of real-world dataset likely represent settings where  $\sigma < 1$  and  $\rho > \rho_0$ . As you pointed out, these settings may favor ECC-S over ECC-Q.*

*We have modified the discussion in Section 5 to provide more detailed suggestions and conclusions regarding the overall performance of the methods.*

## Further changes

- *An error in the specification of the mean vector of the observation in Setting 4 (now Setting 3A) has been corrected: We had set  $\epsilon_0$  to 0, but not specifically mentioned this.*
- *Figure 1 summarizing the univariate results for Settings 1, 2 and 4 (now 1, 3A and 3B) has been updated to reflect the changes in the numbers of the considered settings.*
- *The helpful and constructive comments by all reviewers are now acknowledged.*
- *The bibliography has been updated. In particular, we have added the references Wilks (2015) and Vannitsem et al. (2020). The references Thorarinsdottir and Gneiting (2010) and Hemri and Klein (2017) were moved from the paper to the Supplemental Material as they only occurred in context with Setting 2 (now Setting S1 in the Supplemental Material). Further, we have updated the information for Lang et al. (2019) which is now Lang et al. (2020).*
- *Besides the changes mentioned above, we have made minor changes to various formulations.*

## Changes to the Supplemental Material

*In light of the modifications made to the paper described above, the Supplemental Material has been changed as follows.*

- *A table of contents was added.*
- *Results for additional numbers of members (Settings 1 and 3, now Settings 1 and 2) and dimensions (Setting 1) were added.*
- *Settings 3 and 4 were re-named to Settings 2 and 3A, respectively.*
- *Additional simulation results for the new Setting 3B were added.*
- *Description and results of the simulation setting based on a multivariate truncated Gaussian distribution (Setting 2 in the original submission) were moved to the Supplemental Material where they are accompanied by simulation results for additional parameter choices that were present in the Supplemental Material of the original submission.*

# Simulation-based comparison of multivariate ensemble post-processing methods

Sebastian Lerch<sup>1</sup>, Sándor Baran<sup>2</sup>, Annette Möller<sup>3</sup>, Jürgen Groß<sup>4</sup>, Roman Schefzik<sup>5</sup>, Stephan Hemri<sup>6</sup>, and Maximiliane Graeter<sup>1</sup>

<sup>1</sup>Karlsruhe Institute of Technology, Karlsruhe, Germany

<sup>2</sup>University of Debrecen, Debrecen, Hungary

<sup>3</sup>Technical University of Clausthal, Clausthal, Germany

<sup>4</sup>University of Hildesheim, Hildesheim, Germany

<sup>5</sup>German Cancer Research Center (DKFZ), Heidelberg, Germany

<sup>6</sup>Federal Office of Meteorology and Climatology MeteoSwiss, Zurich-Airport, Switzerland

**Correspondence:** Sebastian Lerch (Sebastian.Lerch@kit.edu)

**Abstract.** Many practical applications of statistical post-processing methods for ensemble weather forecasts require to accurately model spatial, temporal and inter-variable dependencies. Over the past years, a variety of approaches has been proposed to address this need. We provide a comprehensive review and comparison of state of the art methods for multivariate ensemble post-processing. We focus on generally applicable two-step approaches where ensemble predictions are first post-processed separately in each margin, and multivariate dependencies are restored via copula functions in a second step. The comparisons are based on simulation studies tailored to mimic challenges occurring in practical applications and allow to readily interpret the effects of different types of misspecifications in the mean, variance and covariance structure of the ensemble forecasts on the performance of the post-processing methods. Overall, we find that the Schaake shuffle provides a compelling benchmark that is difficult to outperform, whereas the forecast quality of parametric copula approaches and variants of ensemble copula coupling strongly depend on the misspecifications at hand.

*Copyright statement.* TEXT

## 1 Introduction

Despite continued improvements ensemble weather forecasts often exhibit systematic errors that require correction via statistical post-processing methods. Such calibration approaches have been developed for a wealth of weather variables and specific applications. The employed statistical techniques include parametric distributional regression models (Gneiting et al., 2005; Raftery et al., 2005) as well as nonparametric approaches (Taillardat et al., 2016) and semi-parametric methods based on modern machine learning techniques (Rasp and Lerch, 2018). We refer to ~~Vannitsem et al. (2018)~~ [Vannitsem et al. \(2018, 2020\)](#) for a general overview and review.

20 While much of the developments have been focused on univariate methods, many practical applications require to accurately capture spatial, temporal or inter-variable dependencies (Scheffzik et al., 2013). Important examples include hydrological applications (Scheuerer et al., 2017), air traffic management (Chaloulos and Lygeros, 2007) and energy forecasting (Pinson and Messner, 2018). Such dependencies are present in the raw ensemble predictions, but are lost if standard univariate post-processing methods are applied separately in each margin.

25 Over the past years, a variety of multivariate post-processing methods has been proposed, see Scheffzik and Möller (2018) for a recent overview. Those can roughly be categorized into two groups of approaches. The first strategy aims to directly model the joint distribution by fitting a specific multivariate probability distribution. This approach is mostly used in low-dimensional settings, or if a specific structure can be chosen for the application at hand. Examples include multivariate models for temperatures across space (Feldmann et al., 2015), for wind vectors (Schuhen et al., 2012; Lang et al., 2019), and joint  
30 models for temperature and wind speed (Baran and Möller, 2015, 2017).

The second group of approaches proceeds in a two-step strategy. In a first step, univariate post-processing methods are applied independently in all dimensions, and samples are generated from the obtained probability distributions. In a second step, the multivariate dependencies are restored by re-arranging the univariate sample values with respect to the rank order structure of a specific multivariate dependence template. Mathematically, this corresponds to the application of a (parametric  
35 or non-parametric) copula. Examples include ensemble copula coupling (Scheffzik et al., 2013), the Schaake shuffle (Clark et al., 2004) and the Gaussian copula approach (Möller et al., 2013).<sup>1</sup>

Here, we focus on this second strategy which is more generally applicable in cases where no specific assumptions on the parametric structure can be made, or where the dimensionality of the forecasting problem is too high to be handled by fully parametric methods. The overarching goal of this paper is to provide a systematic comparison of state of the art methods for  
40 multivariate ensemble post-processing. In particular, our comparative evaluation includes recently proposed extensions of the popular ensemble copula coupling approach (Hu et al., 2016; Ben Bouallègue et al., 2016). We propose ~~four~~three simulation settings which are tailored to mimic different situations and challenges that arise in applications of post-processing methods. In contrast to case studies based on real-world datasets, simulation studies allow one to specifically tailor the multivariate properties of the ensemble forecasts and observations, and to readily interpret the effects of different types of misspecifications  
45 on the forecast performance of the various post-processing methods. Simulation studies have been frequently applied to analyze model properties, and to compare modeling approaches and verification tools in the context of statistical post-processing, see, e.g., Williams et al. (2014); Thorarinsdottir et al. (2016); Wilks (2017); Allen et al. (2019).

The remainder is organized as follows. Univariate and multivariate post-processing methods are introduced in Section 2. Section 3 provides descriptions of the ~~four~~three simulation settings, with results discussed in Section 4. The paper closes with a  
50 discussion in Section 5. Technical details on specific probability distributions and multivariate evaluation methods are deferred

---

<sup>1</sup>An alternative post-processing approach that allows to preserve multivariate dependencies is the member-by-member method proposed by Van Schaeybroeck and Vannitsem (2015). Scheffzik (2017) demonstrates that member-by-member post-processing can be interpreted as a specific variant of ensemble copula coupling, and can thus be seen as belonging to this group of methods.



to the Appendix. ~~Figures with additional~~ Additional results are available in the Supplementary Material. R (R Core Team, 2019) code with replication material and implementations of all methods is available from [https://github.com/slerch/multiv\\_pp](https://github.com/slerch/multiv_pp).

## 2 Post-processing of ensemble forecasts

We focus on multivariate ensemble post-processing approaches which are based on a combination of univariate post-processing models with copulas. The general two-step strategy of these methods is to first apply univariate post-processing to the ensemble forecasts for each margin (i.e., weather variable, location, and prediction horizon) separately. Then, in a second step, a suitably chosen copula is applied to the univariately post-processed forecasts in order to obtain the desired multivariate post-processing taking account of dependence patterns.

A copula is a multivariate cumulative distribution function (CDF) with standard uniform univariate marginal distributions (Nelsen, 2006). The underlying theoretical background of the above procedure is given by Sklar’s theorem (Sklar, 1959), which states that a multivariate CDF  $H$  (this is what we desire) can be decomposed into a copula function  $C$  modeling the dependence structures (this is what needs to be specified) and its marginal univariate CDFs  $F_1, \dots, F_d$  (this is what is obtained by the univariate post-processing) as follows:

$$H(x_1, \dots, x_d) = C(F_1(x_1), \dots, F_d(x_d))$$

for  $x_1, \dots, x_d \in \mathbb{R}$ . In the approaches considered here, the copula  $C$  is chosen to be either the non-parametric empirical copula induced by a pre-specified dependence template (in the ensemble copula coupling method and variants thereof as well as in the Schaake shuffle), or the parametric Gaussian copula (in the Gaussian copula approach, ~~GCA~~). A Gaussian copula is a particularly convenient parametric model, as apart from the marginal distributions it only requires estimation of the correlation matrix of the multivariate distribution. Under a Gaussian copula the multivariate CDF  $H$  takes the form

$$H(x_1, \dots, x_d | \Sigma) = \Phi_d(\Phi^{-1}(F_1(x_1)), \dots, \Phi^{-1}(F_d(x_d)) | \Sigma), \quad (1)$$

with  $\Phi_d(\cdot | \Sigma)$  denoting the CDF of a  $d$ -dimensional ~~standard~~-normal distribution with mean zero and correlation matrix  $\Sigma$ , and  $\Phi^{-1}$  denoting the quantile function of the univariate standard normal distribution.

To describe the considered methods in more detail in what follows, let  $\mathbf{X}_1, \dots, \mathbf{X}_m \in \mathbb{R}^d$  denote unprocessed ensemble forecasts from  $m$  members, where  $\mathbf{X}_i := (X_i^{(1)}, \dots, X_i^{(d)})$  for  $i = 1, \dots, m$ , and let  $\mathbf{y} := (y^{(1)}, \dots, y^{(d)}) \in \mathbb{R}^d$  be the corresponding verifying observation. We will use  $l = 1, \dots, d$  to denote a multi-index ~~summarizing that may summarize~~ a fixed weather variable, location, and prediction horizon in practical applications to real-world datasets.

### 2.1 Step 1: Univariate post-processing

In a first step, univariate post-processing methods are applied to each margin  $l = 1, \dots, d$  separately. Prominent state-of-the-art univariate post-processing approaches include Bayesian model averaging (Raftery et al., 2005) and ensemble model output statistics (EMOS; Gneiting et al., 2005). In the EMOS approach, which is employed throughout this paper, a non-homogeneous

distributional regression model

$$y^{(l)} | X_1^{(l)}, \dots, X_m^{(l)} \sim F_\theta^{(l)}(y^{(l)} | \theta^{(l)})$$

is fitted, where  $F_\theta^{(l)}$  is a suitably chosen parametric distribution with parameters  $\theta^{(l)} := g(X_1^{(l)}, \dots, X_m^{(l)})$  that depend on the unprocessed ensemble forecast through a link function  $g(\cdot)$ .

85 The choice of  $F_\theta^{(l)}$  is in practice mainly determined by the weather variable being considered in the margin  $l$ . For instance, when  $F_\theta^{(l)}$  can be assumed to be Gaussian with mean  $\mu$  and variance  $\sigma^2$ , such as for temperature or pressure, one may set

$$F_\theta^{(l)} = \mathcal{N}(\mu, \sigma^2), \quad \text{where } (\mu, \sigma^2) := (a_0 + a_1 \bar{X}, b_0 + b_1 S^2) = g(X_1^{(l)}, \dots, X_m^{(l)}) \quad (2)$$

if the ensemble members are exchangeable, with  $\bar{X}$  and  $S^2$  denoting the empirical mean and variance of the ensemble predictions  $X_1^{(l)}, \dots, X_m^{(l)}$ , respectively. The coefficients  $a_0, a_1, b_0$  and  $b_1$  are then derived via suitable estimation techniques using  
90 training data consisting of past ensemble forecasts and observations (Gneiting et al., 2005).

## 2.2 Step 2: Incorporating dependence structures using copulas to obtain multivariate post-processing

When applying univariate post-processing for each margin separately, multivariate (i.e. inter-variable, spatial and/or temporal) dependencies across the margins are lost. These dependencies are restored in a second step. Here, we consider five different approaches to do so. [An overview of selected key features is provided in Table 1. For further discussion of advantages and shortcomings, as well as comparisons of subsets of these methods see, e.g., Schefzik et al. \(2013\); Wilks \(2015\). In the following we use  \$z\$  to denote univariate quantities in the individual dimensions.  \$\mathbf{Z}\$  in bold print is used to represent vector-valued quantities, and  \$Z\$  in normal print is used to for components thereof.](#)  
95

### Assumption of independence (EMOS-Q)

Instead of modeling the desired dependencies in any way, omitting the second step corresponds to assuming independence across the margins. To that end, a univariate sample  $\hat{x}_1^{(l)}, \dots, \hat{x}_m^{(l)}$  is generated in each margin by drawing from the post-processed forecast distribution  $F_\theta^{(l)}, l = 1, \dots, d$ . The univariate samples are then simply combined into a corresponding vector. Following Schefzik et al. (2013), we use equidistant quantiles of  $F_\theta^{(l)}$  at levels  $\frac{1}{m+1}, \dots, \frac{m}{m+1}$  to generate the sample, and denote this approach by EMOS-Q.  
100

### Ensemble copula coupling (ECC)

105 The basic ensemble copula coupling (ECC) approach proposed by Schefzik et al. (2013) proceeds as follows:

1. A sample  $\hat{x}_1^{(l)}, \dots, \hat{x}_m^{(l)}$ , where we assume  $\hat{x}_1^{(l)} \leq \dots \leq \hat{x}_m^{(l)}$  to simplify notation, of the same size  $m$  as the unprocessed ensemble is drawn from each post-processed predictive marginal distribution  $F_\theta^{(l)}, l = 1, \dots, d$ .
2. The sampled values are rearranged in the rank order structure of the raw ensemble, i.e., the permutation  $\sigma_l$  of the set  $\{1, \dots, m\}$  defined by  $\sigma_l(i) = \text{rank}(X_i^{(l)})$ , with possible ties resolved at random, is applied to the post-processed sample

**Table 1.** Overview of selected key characteristics of the multivariate post-processing methods considered in this paper.

<u>Method</u>	<u>Dependence template</u>	<u>Flow-dependent copula structure</u>	<u>Size of resulting multivariate ensemble</u>	<u>Univariate sampling</u>	<u>Involves randomness</u>
<u>EMOS-Q</u>	<u>assumes independence</u>	<u>~</u>	<u>arbitrary</u>	<u>equidistant</u>	<u>no</u>
<u>ECC-R</u>	<u>raw ensemble</u>	<u>yes</u>	<u><math>m</math></u>	<u>random</u>	<u>yes (sampling)</u>
<u>ECC-Q</u>	<u>raw ensemble</u>	<u>yes</u>	<u><math>m</math></u>	<u>equidistant</u>	<u>no</u>
<u>ECC-S</u>	<u>raw ensemble</u>	<u>yes</u>	<u><math>m</math></u>	<u>stratified</u>	<u>yes (sampling)</u>
<u>dECC</u>	<u>raw ensemble &amp; forecast errors</u>	<u>yes</u>	<u><math>m</math></u>	<u>equidistant</u>	<u>no</u>
<u>SSh</u>	<u>observations</u>	<u>no</u>	<u>arbitrary</u>	<u>equidistant</u>	<u>yes (selection of training cases)</u>
<u>GCA</u>	<u>observations</u>	<u>no</u>	<u>arbitrary</u>	<u>random</u>	<u>yes (sampling)</u>

110 from the first step in order to obtain the final ECC ensemble  $\tilde{X}_1^{(l)}, \dots, \tilde{X}_m^{(l)}$  via

$$\tilde{X}_i^{(l)} = \hat{x}_{\sigma_l(i)}^{(l)},$$

where  $i = 1, \dots, m$  and  $l = 1, \dots, d$ .

Depending on the specific sampling procedure in step 1, we here distinguish the following different ECC variants:

- **ECC-R:** The sample  $\hat{x}_1^{(l)}, \dots, \hat{x}_m^{(l)}$  is randomly drawn from  $F_\theta^{(l)}$  (and subsequently arranged in ascending order).
- 115 – **ECC-Q:** The sample is constructed using equidistant quantiles of  $F_\theta^{(l)}$  at levels  $\frac{1}{m+1}, \dots, \frac{m}{m+1}$ :

$$\hat{x}_1^{(l)} := (F_\theta^{(l)})^{-1}\left(\frac{1}{m+1}\right), \dots, \hat{x}_m^{(l)} := (F_\theta^{(l)})^{-1}\left(\frac{m}{m+1}\right).$$

- **ECC-S** (Hu et al., 2016): First, random numbers  $u_1, \dots, u_m$ , where  $u_i \sim \mathcal{U}\left(\frac{i-1}{m}, \frac{i}{m}\right)$  for  $i = 1, \dots, m$ , are drawn, with  $\mathcal{U}(a, b]$  denoting the uniform distribution on the interval  $(a, b]$ . Then,  $\hat{x}_i^{(l)}$  is set to the quantile of  $F_\theta^{(l)}$  at level  $u_i$ :

$$\hat{x}_1^{(l)} := (F_\theta^{(l)})^{-1}(u_1), \dots, \hat{x}_m^{(l)} := (F_\theta^{(l)})^{-1}(u_m).$$

120 Besides the above sampling schemes, Scheffzik et al. (2013) propose an alternative transformation approach referred to as ECC-T. This variant is in particular appealing for theoretical considerations, as it provides a link between the ECC notion and member-by-member post-processing approaches (Scheffzik, 2017). However, as it may involve additional modeling steps, ECC-T is not as generic as the other schemes and thus not explicitly considered [in this paper here](#).

### Dual ensemble copula coupling (dECC)

125 Dual ECC (dECC) is an extension of ECC which aims at combining the structure of the unprocessed ensemble with a component accounting for the forecast error autocorrelation structure (Ben Bouallègue et al., 2016), proceeding as follows:

1. ECC-Q is applied in order to obtain re-ordered ensemble forecasts  $\tilde{\mathbf{X}}_1, \dots, \tilde{\mathbf{X}}_m$ , with  $\tilde{\mathbf{X}}_i := (\tilde{X}_i^{(1)}, \dots, \tilde{X}_i^{(d)})$  for  $i = 1, \dots, m$ .
2. A transformation based on an estimate of the error autocorrelation  $\hat{\Sigma}_e$  is applied to the bias-corrected post-processed forecast in order to obtain correction terms  $\mathbf{c}_1, \dots, \mathbf{c}_m$ . Precisely,  $\mathbf{c}_i := (\hat{\Sigma}_e)^{\frac{1}{2}} \cdot (\tilde{\mathbf{X}}_i - \mathbf{X}_i)$  for  $i = 1, \dots, m$ .
- 130 3. An adjusted ensemble  $\check{\mathbf{X}}_1, \dots, \check{\mathbf{X}}_m$  is derived via  $\check{\mathbf{X}}_i := \tilde{\mathbf{X}}_i + \mathbf{c}_i$  for  $i = 1, \dots, m$ .
4. ECC-Q is applied again, but now performing the re-ordering with respect to the rank order structure of the adjusted ensemble from step 3 used as a modified dependence template.

### Schaake shuffle (SSh)

135 The Schaake shuffle (SSh) proceeds as ECC-Q, but re-orders the sampled values in the rank order structure of  $m$  ~~randomly determined~~ past observations (Clark et al., 2004) and not with respect to the unprocessed ensemble forecasts. For a better comparison with (d)ECC, the size of the SSh ensemble is restricted to equal that of the unprocessed ensemble here. However, in principle, the SSh ensemble may have an arbitrary size, provided that sufficiently [enough many](#) past observations are available to build the dependence template. Extensions of the SSh that select past observations based on similarity are available (Scheffzik, 140 2016; Scheuerer et al., 2017), but not explicitly considered here as their implementation is not ~~directly~~ straightforward and may involve additional modelling choices specific to the situation at hand.

The reordering-based methods considered thus far can be interpreted as non-parametric, empirical copula approaches. In particular, in the setting of Sklar's theorem,  $C$  is taken to be the empirical copula induced by the corresponding dependence template, i.e., the unprocessed ensemble forecasts in case of ECC, the adjusted ensemble in case of dECC, and the past 145 observations in case of the SSh.

### Gaussian copula approach (GCA)

By contrast, in the Gaussian copula approach (GCA) proposed by Pinson and Girard (2012) and Möller et al. (2013), the copula  $C$  is taken to be the parametric Gaussian copula. GCA can be traced back to similar ideas from earlier work in spatial statistics (e.g., Berrocal et al., 2008) and proceeds as follows:

150 1. A set of past observations  $\mathbf{y}_1, \dots, \mathbf{y}_K$ , with  $\mathbf{y}_k = (y_k^{(1)}, \dots, y_k^{(d)})$ , is transformed into latent standard Gaussian observations  $\tilde{\mathbf{y}}_1, \dots, \tilde{\mathbf{y}}_K$  by setting

$$\tilde{y}_k^{(l)} = \Phi^{-1} \left( F_\theta^{(l)}(y_k^{(l)}) \right) \quad (3)$$

for  $k = 1, \dots, K$  and  $l = 1, \dots, d$ , where  ~~$\Phi^{-1}$  is the inverse of the CDF of the univariate standard normal distribution and~~  $F_\theta^{(l)}$  is the marginal distribution obtained by univariate post-processing. The index  $k = 1, \dots, K$  here refers to a training  
155 set of past observations.

2. An empirical (or parametric)  $(d \times d)$  correlation matrix  $\widehat{\Sigma}$  of the  $d$ -dimensional ~~standard~~ normal distribution in (1) is estimated from  $\tilde{\mathbf{y}}_1, \dots, \tilde{\mathbf{y}}_K$ .

3. Multivariate random samples  $\mathbf{Z}_1, \dots, \mathbf{Z}_m \sim \mathcal{N}_d(\mathbf{0}, \widehat{\Sigma})$  are drawn, where  $\mathcal{N}_d(\mathbf{0}, \widehat{\Sigma})$  denotes a  $d$ -dimensional normal distribution with mean vector  $\mathbf{0} := (0, \dots, 0)$  and estimated correlation matrix  $\widehat{\Sigma}$  from Step 2, and  $\mathbf{Z}_i := (Z_i^{(1)}, \dots, Z_i^{(d)})$   
160 for  $i = 1, \dots, m$ .

4. The final GCA post-processed ensemble forecast  $\mathbf{X}_1^*, \dots, \mathbf{X}_m^*$ , with  $\mathbf{X}_i^* := (X_i^{*(1)}, \dots, X_i^{*(d)})$  for  $i = 1, \dots, m$  is obtained via

$$X_i^{*(l)} := \left( F_\theta^{(l)} \right)^{-1} \left( \Phi(Z_i^{(l)}) \right) \quad (4)$$

for  $i = 1, \dots, m$  and  $l = 1, \dots, d$ , with  $\Phi$  denoting the CDF of the univariate standard normal distribution. While the size  
165 of the resulting ensemble may in principle be arbitrary, it is here set to the size  $m$  of the raw ensemble.

### 3 Simulation settings

We consider several simulation settings to highlight different aspects and provide a broad comparison of the effects of potential misspecifications of the ensemble predictions on the performance of the various multivariate post-processing methods. The general setup of all simulation settings is as follows.

170 An initial training set of pairs of simulated ensemble forecasts and observations of size  $n_{\text{init}}$  is generated. Post-processed forecasts are then computed and evaluated over a test set of size  $n_{\text{test}}$ . Therefore,  $n := n_{\text{init}} + n_{\text{test}}$  iterations are performed in total for all simulation settings. In the following, we set  $m = 50, n_{\text{init}} = 500, n_{\text{test}} = 1000$  throughout.

To describe the individual settings in more detail, we here begin by first identifying the general structure of the steps that are performed in all settings. For each iteration  $t$  in both training and test set (i.e.,  $t = 1, \dots, n$ ), multivariate forecasts and  
175 observations are generated:

(S1) Generate multivariate observations and ensemble forecasts.

For all iterations  $t$  in the test set (i.e.,  $t = n_{\text{init}} + 1, \dots, n$ ), the following steps are carried out:

(S2) Apply univariate post-processing separately in each dimension.<sup>2</sup>

(S3) Apply multivariate post-processing methods.

180 (S4) Compute univariate and multivariate measures of forecast performance on the test set.

Unless indicated otherwise all simulation draws are independent across iterations. To simplify notation we will thus typically omit the simulation iteration index  $t$  in the following.

To quantify simulation uncertainty, the above procedure is repeated 100 times for each tuning parameter combination in each setting. In the interest of brevity, we omit ECC-R which did show substantially worse results in initial tests (see also Schefzik  
185 et al., 2013). In the following, the individual simulation settings are described in detail, and specific implementation choices are discussed.

### 3.1 Setting 1: Multivariate Gaussian distribution

As starting point we first consider a simulation model where observations and ensemble forecasts are drawn from multivariate Gaussian distributions.<sup>3</sup> [This setting may for example apply in the case of temperature forecasts at multiple locations considered](#)  
190 [simultaneously](#). The simplicity of this model allows to readily interpret misspecifications in the mean, variance and covariance structures.

(S1) For iterations  $t = 1, \dots, n$ , independent and identically distributed samples of observations and ensemble forecasts are generated as follows:

- observation:  $\mathbf{y} \sim \mathcal{N}_d(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}^0)$ , where  $\boldsymbol{\mu}_0 = (0, \dots, 0) \in \mathbb{R}^d$ , and  $\Sigma_{i,j}^0 = \rho_0^{|i-j|}$ , for  $i, j = 1, \dots, d$ .
- 195 – ensemble forecasts:  $\mathbf{X}_1, \dots, \mathbf{X}_m \stackrel{\text{iid}}{\sim} \mathcal{N}_d(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , where  $\boldsymbol{\mu} = (\epsilon, \dots, \epsilon) \in \mathbb{R}^d$ , and  $\Sigma_{i,j} = \sigma \rho^{|i-j|}$ , for  $i, j = 1, \dots, d$ .

The parameters  $\epsilon$  and  $\sigma$  introduce a bias and a misspecified variance in the marginal distributions of the ensemble forecasts. These systematic errors are kept constant across dimensions  $1, \dots, d$ . The parameters  $\rho_0$  and  $\rho$  control the autoregressive structure of the correlation matrix of the observations and ensemble forecasts. Setting  $\rho_0 \neq \rho$  introduces misspecifications of the correlation structure of the ensemble forecasts.

200 (S2) As described in Section 2.1, univariate post-processing is applied independently in each dimension  $1, \dots, d$ . Here, we employ the standard Gaussian EMOS model (2) proposed by Gneiting et al. (2005). The EMOS coefficients  $a_0, a_1, b_0, b_1$  are estimated by minimizing the mean continuous ranked probability score (CRPS, see Appendix B) over the training set consisting of the  $n_{\text{init}}$  initial iterations, and are then used to produce out of sample forecasts for the  $n_{\text{test}}$  iterations in the test set.

---

<sup>2</sup>With the exception of Setting 4.3, the estimation of univariate post-processing models utilizes the initial training set only. Setting 4.3 covers the possibly more realistic case of variations across repetitions of the experiment.

<sup>3</sup>Wilks (2017) considers a similar setting in the context of multivariate calibration assessment which we here extend towards multivariate ensemble post-processing.

205 (S3) Next, the multivariate post-processing methods described in Section 2.2 are applied. Implementation details for the individual methods are as follows.

- For dECC, the estimate of the error autocorrelation  $\hat{\Sigma}_e$  is obtained from the  $n_{\text{init}}$  initial training iterations to compute the required correction terms for the test set.
- To obtain the dependence template for SSh,  $m$  past observations are randomly selected from all iterations preceding the current iteration  $t$ .
- The correlation matrix  $\Sigma$  required for GCA is estimated by the empirical correlation matrix based on all iterations preceding the current iteration  $t$ .
- The verification results for all methods that require random sampling (ECC-S, SSh, GCA) are averaged over 10 independent repetitions for each iteration  $t = n_{\text{init}} + 1, \dots, n$  in the test set.

215 The multivariate Gaussian setting is implemented for  $d = 5$  and all combinations of  $\epsilon \in \{0, 1, 3\}$ ,  $\sigma^2 \in \{0.5, 1, 2, 5\}$ , and  $\rho, \rho_0 \in \{0.1, 0.25, 0.5, 0.75, 0.9\}$ . As indicated above, the simulation experiment is repeated 100 times for each of the 300 parameter combinations.

### 3.2 ~~Setting 2: Multivariate truncated Gaussian distribution~~

220 ~~Setting 1 can be generalized by replacing the multivariate Gaussian distribution by a multivariate truncated Gaussian distribution  $\mathcal{N}_{a,b}^d(\mu, \Sigma)$ , where  $a$  and  $b$  are the vectors of lower and upper truncation points, respectively. In univariate settings this distribution plays important role in wind speed modelling (Thorarinsdottir and Gneiting, 2010) or in [If the setting from above is interpreted as a multivariate model for temperatures at multiple locations, observations from the extant literature on post-processing of hydrological forecasts \(Hemri and Klein, 2017\)](#). Compared to Setting 1, here misspecifications in location vector  $\mu$  and /or scale matrix  $\Sigma$  result in more complex deviations in mean vectors and covariance matrices. [suggest that typically, values of  \$\sigma < 1\$  and  \$\rho > \rho\_0\$  would be expected in real-world datasets.](#)~~

225 ~~For iterations  $t = 1, \dots, n$ , independent and identically distributed samples of observations  $\mathbf{y}$  and ensemble forecasts  $\mathbf{X}_1, \dots, \mathbf{X}_m$  are generated from  $\mathcal{N}_{a,b}^d(\mu_0, \Sigma^0)$  and  $\mathcal{N}_{a,b}^d(\mu, \Sigma)$ , respectively, where  $\Sigma^0$  and  $\Sigma$  are defined as in Setting 1,  $\mu_0 = (\mu_0, \dots, \mu_0) \in \mathbb{R}^d$  and  $\mu = (\mu, \dots, \mu) \in \mathbb{R}^d$ . Univariate post-processing is based on the truncated normal EMOS model of Hemri and Klein (2017), where the EMOS coefficients are calculated by optimizing the mean CRPS over the training set consisting of the  $n_{\text{init}}$  initial iterations. Similar to~~

230 ~~[A variant of Setting 1](#), the obtained EMOS models are used to produce out-of sample forecasts for the  $n_{\text{test}}$  iterations in the test set. [Identical to \(S3\) of Setting based on a multivariate truncated Gaussian distribution has also been investigated. Apart from a slightly worse performance of GCA, the results are similar to those of Setting 1. For simplicity, we consider a lower truncation at 0 only, i.e.,  \$a = \(0, \dots, 0\)\$](#)  We thus refer to Section 5 of the Supplemental Material where details on the simulation [setting and  \$b = \(-\infty, \dots, \infty\)\$](#) . The truncated Gaussian setting is implemented for  $d = 5$  and all combinations of~~

~~$\mu_0 \in \{2, 3\}$ ,  $\mu \in \{2, 3, 5\}$ ,  $\rho_0 \in \{0.25, 0.5, 0.75\}$ ,  $\rho \in \{0.1, 0.25, 0.5, 0.75, 0.9\}$ ,  $\sigma \in \{0.25, 0.5, 1, 3, 5\}$ ,~~

resulting in 450 experiments which are repeated 100 times each results are provided.

### 3.2 Setting 32: Multivariate censored extreme value distribution

To investigate alternative marginal distributions employed in post-processing applications, we further consider a simulation setting based on a censored version of the generalized extreme value (GEV) distribution. The GEV distribution was introduced by Jenkinson (1955) among others, combining three different types of extreme value distributions. It has been widely used for modelling extremal climatological events such as flood peaks (e.g., Morrison and Smith, 2002) or extreme precipitation (e.g., Feng et al., 2007). In the context of post-processing, GEV distributions have for example been applied for modeling wind speed in Lerch and Thorarinsdottir (2013). Here, we consider multivariate observations and forecasts with marginal distributions given by a left-censored version of the GEV distribution which was proposed by Scheuerer (2014) in the context of post-processing ensemble forecasts of precipitation amounts.

(S1) For iterations  $t = 1, \dots, n$  samples of observations and ensemble forecasts are generated as follows. For  $l = 1, \dots, d$ , the marginal distributions are GEV distributions left-censored at 0,

$$F_{\theta}^{(l)} = \text{GEV}_0(\mu, \sigma, \xi),$$

where the distribution parameters  $\mu$  (location),  $\sigma$  (scale) and  $\xi$  (shape) are identical across dimensions  $l = 1, \dots, d$ . Details on the left-censored GEV distribution are provided in Appendix A. Misspecifications of the marginal ensemble predictions are obtained by choosing different GEV parameters for observations  ~~$(\mu_y, \sigma_y, \xi_y)$~~  and forecasts  ~~$(\mu_x, \sigma_x, \xi_x)$~~   $(\mu_0, \sigma_0, \xi_0)$  and forecasts  $(\mu, \sigma, \xi)$ . Combined misspecifications of the three parameters result in more complex deviations of mean and variance (on the univariate level) especially compared to Setting 1, but also compared to Setting 2-1. Typically there is a joint influence of the GEV parameters on mean and dispersion properties of the distribution. In order to exploit the complex behavior a variety of parameter combinations for observations and ensemble forecasts were considered.

To generate multivariate observations  $\mathbf{y} = (y^{(1)}, \dots, y^{(d)})$  and ensemble predictions  $\mathbf{X}_i = (X_i^{(1)}, \dots, X_i^{(d)})$ ,  $i = 1, \dots, m$ , the so-called NORTA (normal to anything) approach is chosen, see Cario and Nelson (1997); Chen (2001). This method allows one to generate realizations of a random vector  $\mathbf{z} = (z^{(1)}, \dots, z^{(d)})$  with specified marginal distribution functions  $F_{\theta}^{(l)}$ ,  $l = 1, \dots, d$ , and a given correlation matrix  ~~$\mathbf{R} = (\text{Corr}(z^{(k)}, z^{(l)}))_{k,l=1}^d$~~ .

$\mathbf{R} = (\text{Corr}(z^{(k)}, z^{(l)}))_{k,l=1}^d$ . The NORTA procedure consists of three steps. In a first step a vector  $\mathbf{v} = (v^{(1)}, \dots, v^{(d)})$  is generated from  ~~$\mathcal{N}_d(0, \mathbf{R}^*)$~~   $\mathcal{N}_d(0, \mathbf{R}^*)$  for a correlation matrix  ~~$\mathbf{R}^*$~~   $\mathbf{R}^*$ . In a second step,  $u^{(l)} = \Phi(v^{(l)})$  is computed, where  $\Phi$  denotes the CDF of the standard normal distribution. In a third step,  $z^{(l)} = \left(F_{\theta}^{(l)}\right)^{-1}(u^{(l)})$  is derived for  $l = 1, \dots, d$ , where  $\left(F_{\theta}^{(l)}\right)^{-1}$  is the inverse of  $F_{\theta}^{(l)}$ . The correlation matrix  ~~$\mathbf{R}^*$~~   $\mathbf{R}^*$  is chosen in a such a way that the  $z^{(l)}$  have the desired target correlation matrix  $\mathbf{R}$ . Naturally, the specification of  ~~$\mathbf{R}^*$~~   $\mathbf{R}^*$  is the most involved part of this procedure. Here, we use the retrospective approximation algorithm implemented in the R package NORTARA (Su, 2014). The NORTARA package infrequently produced error and warnings, which were not present for alternative starting values of the random number generator.



Following the previous simulation settings the target correlation matrix  $R$  is chosen as

270

$$R_{i,j} = \rho^{|i-j|}$$

$$R_{i,j} = \rho^{|i-j|} \text{ for } -1 < \rho < 1 \text{ and } i, j = 1, \dots, d.$$

- (S2) To separately post-process the univariate ensemble forecasts we employ the EMOS method for quantitative precipitation based on the left-censored GEV distribution proposed by Scheuerer (2014). To that end we assume  $-0.278 < \xi < 0.5$ , in which case such that the mean  $\nu$  and the variance of the non-censored GEV distribution exist, and

275

$$\nu = \begin{cases} \mu + \sigma \frac{\Gamma(1-\xi)-1}{\xi}, & \xi \neq 0 \\ \mu + \sigma\gamma, & \xi = 0 \end{cases},$$

where  $\Gamma$  denotes the gamma function and  $\gamma$  is the Euler-Mascheroni constant. See Appendix A for comments on mean and variance of the left-censored GEV. Following Scheuerer (2014), the parameters  $(\nu, \sigma, \xi)$  are linked to the ensemble predictions via

$$g(X_1^{(l)}, \dots, X_m^{(l)}) = \left( a_0 + a_1 \bar{X}^{(l)} + a_2 \bar{X}_{\underline{0}z}^{(l)}, b_0 + b_1 \text{MD}_X^{(l)}, \xi \right).$$

Here,  $\bar{X}^{(l)}$  and  $\bar{X}_{\underline{0}z}^{(l)}$  are the arithmetic mean and the fraction of zero values of the ensemble predictions  $X_1^{(l)}, \dots, X_m^{(l)}$ , respectively, while  $\text{MD}_X^{(l)}$  denotes the mean absolute difference of the ensemble predictions, i.e.,

$$\text{MD}_X^{(l)} = \frac{1}{m^2} \sum_{i=1}^m \sum_{j=1}^m |X_i^{(l)} - X_j^{(l)}|.$$

The shape parameter  $\xi$  is not linked to the ensemble predictions, but is estimated along with the EMOS coefficients  $a_0, a_1, a_2$  and  $b_0, b_1$ . As in Scheuerer (2014), the link function refers to the parameter  $\nu$  instead of  $\mu$ , since it is argued that for fixed  $\nu$  an increase in  $\sigma$  can be interpreted more naturally as an increase in uncertainty.

280

An implementation in R is available in the `ensembleMOS` package (Yuen et al., 2018). For our simulation, this package was not directly invoked, but the respective functions were used as a template. As described in Section 2.1, univariate post-processing is applied independently in each dimension  $l = 1, \dots, d$ . The EMOS coefficients are estimated as described above over the training set consisting of the  $n_{\text{init}}$  initial iterations, and are then used to produce out of sample forecasts for the  $n_{\text{test}}$  iterations in the test set.

285

- (S3) Identical to (S3) of Setting 1, except for GCA, where we proceed differently to account for the point mass at zero. The latent standard Gaussian observations  $\tilde{y}_k^{(l)}$  are generated by  $\tilde{y}_k^{(l)} = \Phi^{-1}(u)$ , where  $u$  is a randomly chosen value in the interval  $(0, F_\theta^{(l)}(0))$  in case  $y_k^{(l)} = 0$  and  $u = F_\theta^{(l)}(y_k^{(l)})$  in case  $y_k^{(l)} > 0$ .

The multivariate censored extreme value setting is implemented for  $d = 4$  and four different scenarios summarized in Table 2. The choice of dimension is motivated by the fact that preliminary analyses had revealed a heavy increase of computation time and numerical problems for values of  $d$  greater than 4. In each scenario the  $\text{GEV}_0$  distribution parameters for the observations

	$\mu_y$	$\xi_y$	$\sigma_y$	$\mu_x$	$\xi_x$	$\sigma_x$
A	0.0	-0.1	1.0	1.0	0.0	0.2
B	0.0	-0.1	1.0	0.0	0.0	2.0
C	1.0	0.3	1.0	0.0	0.0	2.0
D	0.0	0.0	1.0	0.0	0.0	1.0

**Table 2.** Different simulation scenarios for Setting 3.2.

290 are chosen according to  $(\mu_y, \xi_y, \sigma_y)(\mu_0, \xi_0, \sigma_0)$ , while the parameters for the ensemble predictions are chosen according to  $(\mu_x, \xi_x, \sigma_x)(\mu, \xi, \sigma)$ . In both cases, the correlation matrix  $R$  from above is invoked with different choices of  $\rho_y$  and  $\rho_x$  and  $\rho$  from the set  $\{0.25, 0.5, 0.75\}$  giving a total of  $4 \times 9 = 36$  scenarios.

Note that according to Scheuerer (2014) there is a positive probability for zero to occur when either  $\xi \leq 0$  or  $\xi > 0$  and  $\mu < \sigma/\xi$ . The scenarios from Table 2 are chosen in such a way that either one of these two conditions is met.

295 The scenarios from Table 2 were not chosen to mimic real life situations in the first place, but to emulate pronounced differences in distributions and account for a variety of misspecification types. In future research a more detailed and data based study of the properties of the GEV<sub>0</sub> in ensemble postprocessing of precipitation is planned, which might give further insight into the correspondence (and interplay) of the GEV<sub>0</sub> parameters to typically occurring situations for precipitation.

### 3.3 Setting 43: Multivariate Gaussian distribution with changes over iterationtime

300 In the preceding simulation settings, the misspecifications of the ensemble forecasts were kept constant over the iterations  $t = 1, \dots, n$  within the simulation experiments. However, forecast errors of real-world ensemble predictions often exhibit systematic changes over time, for example due to seasonal effects or differences in flow-dependent predictability due to variations of large scale atmospheric conditions. Here, we modify the multivariate Gaussian simulation setting from Section 3.1 to introduce changes in the mean, variance and covariance structure of the multivariate distributions of observations and ensemble forecasts. In analogy to practical applications of multivariate post-processing, the ensemble predictions and observations may be interpreted as multivariate in terms of location or prediction horizon, with changes of the misspecification properties over time.

(S1) For iterations  $t = 1, \dots, n$ , independent samples of observations and ensemble forecasts are generated as follows:

- observation:  $\mathbf{y} \sim \mathcal{N}_d(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}^0)$ , where  $\boldsymbol{\mu}_0 = \sin\left(\frac{2\pi t}{n}\right) + (\epsilon_0, \dots, \epsilon_0)^T \in \mathbb{R}^d$ ,  $\boldsymbol{\mu}_0 = \sin\left(\frac{2\pi t}{n}\right) + (0, \dots, 0)^T \in \mathbb{R}^d$ . To obtain the correlation matrix  $\boldsymbol{\Sigma}^0$ , let  $R_{i,j} = \rho_0^{|i-j|} + \sin\left(\frac{2\pi t}{n}\right)$ , for  $i, j = 1, \dots, d$  and  $\mathbf{S}_0 = \mathbf{R}\mathbf{R}^T$ ,  $\mathbf{S}_0 = \mathbf{R}\mathbf{R}^T$ . The covariance matrix  $\mathbf{S}_0$  is scaled into the corresponding correlation matrix  $\boldsymbol{\Sigma}^0$  using the R function `cov2cor()`.
- ensemble forecasts:  $\mathbf{X}_1, \dots, \mathbf{X}_m \stackrel{\text{iid}}{\sim} \mathcal{N}_d(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , where  $\boldsymbol{\mu} = \sin\left(\frac{2\pi t}{n}\right) + (\epsilon, \dots, \epsilon)^T \in \mathbb{R}^d$ . To obtain the correlation matrix  $\boldsymbol{\Sigma}$  we proceed as for the observations, however, we set  $R_{i,j} = \rho^{|i-j|} + \sin\left(\frac{2\pi t}{n}\right)$ , for  $i, j = 1, \dots, d$  (i.e.,  $\rho_0$  is replaced by  $\rho$ ).

315 In contrast to Setting 1, the misspecifications in the mean and correlation structure now include a periodic component.  
The above setup will be denoted by Setting 3A.

Following a suggestion from an anonymous reviewer, we further consider a variant which we refer to as Setting 3B. For iterations  $t = 1, \dots, n$  we generate independent samples of observations and ensemble forecasts as follows:

- observation:  $\mathbf{y} \sim \mathcal{N}_d(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}^0(t))$ , where  $\boldsymbol{\mu}_0 = (0, \dots, 0)^T \in \mathbb{R}^d$ . To obtain a correlation matrix  $\boldsymbol{\Sigma}(t)$  that varies over iterations we set  $\Sigma_{i,j}^0(t) = \rho_0^{|i-j|}(t)$ , for  $i, j = 1, \dots, d$ , where the correlation parameter  $\rho_0(t)$  varies over iterations according to

320

$$\rho_0(t) = \rho_0 \cdot \left(1 - \frac{a}{2}\right) + \rho_0 \cdot \left(\frac{a}{2}\right) \sin\left(\frac{2\pi t}{n}\right)$$

for  $a \in (0, 1)$ . The lag-1 correlations thus oscillate between  $\rho_0$  and  $\rho_0 \cdot (1 - a)$ .

- ensemble forecasts:  $\mathbf{X}_1, \dots, \mathbf{X}_m \stackrel{\text{iid}}{\sim} \mathcal{N}_d(\boldsymbol{\mu}, \boldsymbol{\Sigma}(t))$ , where  $\boldsymbol{\mu} = (\epsilon, \dots, \epsilon)^T \in \mathbb{R}^d$ . Similar to the observations, we set  $\Sigma_{i,j}(t) = \sigma \rho^{|i-j|}(t)$ , for  $i, j = 1, \dots, d$ , where

325

$$\rho(t) = \rho \cdot \left(1 - \frac{a}{2}\right) + \rho \cdot \left(\frac{a}{2}\right) \sin\left(\frac{2\pi t}{n}\right),$$

with  $a$  from above. The correlations for the ensemble member forecasts thus oscillate between  $\rho$  and  $\rho \cdot (1 - a)$

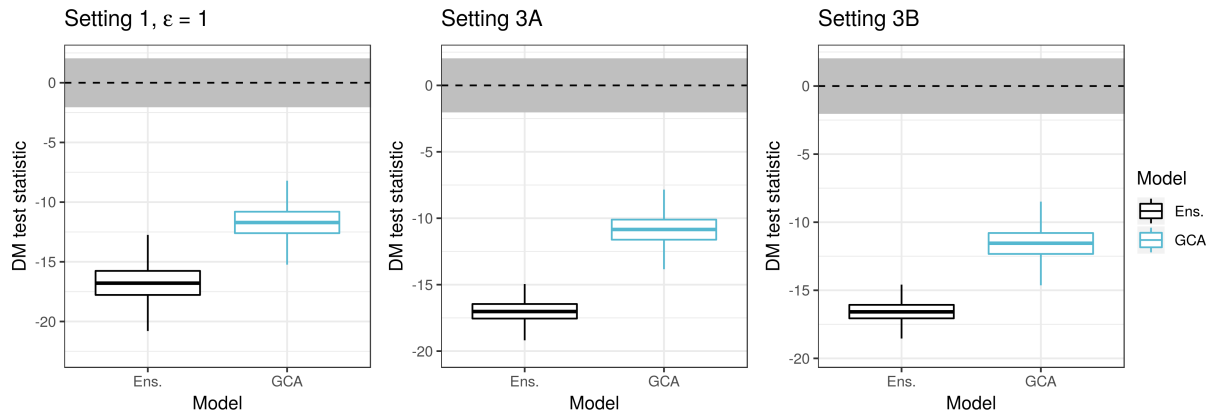
Settings 3A and 3B differ in the variations of the mean and covariance structure over time. For both, we proceed as follows.

330 (S2) As in Setting 1, we employ the standard Gaussian EMOS model (2). However, to account for the changes over iterations we now utilize a rolling window consisting of pairs of ensemble forecasts and observations from the 100 iterations preceding  $t$  as training set to obtain estimates of the EMOS coefficients. See Lang et al. (2020) for a detailed discussion of alternative approaches to incorporate time dependence in the estimation of post-processing models.

335 (S3) The application of the multivariate post-processing methods is identical to the approach taken in Setting 1. Note that we deliberately follow the naive standard implementations (see Section 2.2) here to highlight some potential issues of the Schaake shuffle in this context.

~~The above setting Setting 3A~~ is implemented for  ~~$d = 5, \epsilon = 1, \sigma^2 = 1$~~   $d = 5, \epsilon = 1$  and all combinations of  $\rho, \rho_0 \in \{0.1, 0.25, 0.5, 0.75, 0.9\}$ .  
~~As before the simulation experiment is~~ For Setting 3B, we investigate separate sets of low ( $\rho_0 = 0.25$ ), medium ( $\rho_0 = 0.5$ ) and high ( $\rho_0 = 0.75$ ) true correlation, with corresponding choices of  $\rho$  with low ( $\rho \in \{0.2, 0.25, 0.3\}$ ), medium ( $\rho \in \{0.4, 0.45, 0.5, 0.55, 0.6\}$ ) and high ( $\rho \in \{0.7, 0.75, 0.8\}$ ) values, respectively. Further, values of  $d = 5, \epsilon = 1, \sigma \in \{0.5, 1, 5\}$  and  $a \in \{0.2, 0.5, 0.7\}$  are considered for each of these sets. As before simulation experiments are repeated 100 times for each of the parameter combinations.

340



**Figure 1.** Summaries of DM test statistic values based on the CRPS. ECC-Q forecasts are used as reference model such that positive values of the test statistic indicate improvements over ECC-Q and negative ~~value indicated~~ values indicate deterioration of forecast skill. Boxplots summarize results ~~of~~ from multiple parameter combinations for the simulation settings, with potential restrictions on the simulation parameters indicated in the plot title. For example, boxplots in the first panel summarize simulation results from all parameter combinations of Setting 1 (and the 100 Monte Carlo repetitions each) subject to  $\epsilon = 1$ . The horizontal gray stripe indicates the acceptance region of the two-sided DM test under the null hypothesis of equal predictive performance at a level of 0.05.

## 4 Results

In the following, we focus on comparisons of the relative predictive performance of the different multivariate post-processing methods and apply proper scoring rules for forecast evaluation. In particular, we use the energy score (ES; Gneiting et al., 2008) and variogram score of order 1 (VS; Scheuerer and Hamill, 2015) to evaluate multivariate forecast performance. Diebold-Mariano (DM; Diebold and Mariano, 1995) tests are applied to assess the statistical significance of the score differences between models. Details on forecast evaluation based on proper scoring rules and DM tests are provided in Appendix B. Note that proper scoring rules are often used in the form of skill scores to investigate relative improvements in predictive performance in the meteorological literature. Here, we instead follow suggestions of Ziel and Berk (2019) who argue that the use of DM tests is of crucial importance to appropriately discriminate between multivariate models.

While our focus here is on multivariate performance, we briefly demonstrate that the univariate post-processing models applied in the different simulation settings usually work as intended.

### 4.1 Univariate performance

The univariate predictive performance of the raw ensemble forecasts in terms of the CRPS is improved by the application of univariate post-processing methods across all parameter choices in all simulation settings. The magnitude of the relative improvements by post-processing depends on the chosen simulation parameters, exemplary results are shown in Figure 1. ~~Note that~~ The results for Setting 2 are omitted as they vary more and strongly depend on the simulation parameters.

360 ECC-Q does not change the marginal distributions, the univariate forecasts are thus identical to solely applying univariate  
post-processing methods in the margins separately, without accounting for dependencies. We will later refer to this as EMOS-Q.  
Note that for ECC-S, ~~dECC and SSh are omitted as~~ and SSh differences in the univariate forecast distributions ~~are identical~~  
compared to those of ECC-Q, ~~subject may arise from randomly sampling the quantile levels in ECC-S and due~~ to random  
fluctuations due to the 10 random repetitions that were performed to account for simulation uncertainty of those methods. <sup>4</sup>  
However, we found the effects on the univariate results to be negligible and omit ECC-S, dECC and SSh from Figure 1.

365 For the simulation parameter values summarized there, univariate post-processing works as intended with statistically significant improvements over the raw ensemble forecasts. Note that for GCA the univariate marginal distributions are modified due to the transformation step in (4). While the quantile forecasts of ECC-Q are close to optimal in terms of the CRPS (Bröcker, 2012) the (randomly sampled) univariate GCA forecasts do not possess this property, resulting in worse univariate performance compared to all other methods.

370 ~~To give an impression of the univariate performance of raw ensemble forecasts compared to post-processing for Setting 3, the four scenarios from Table 2 are considered and CRPS skill scores for the raw ensemble are computed with ECC-Q as reference model, see Figure ??.~~ Negative skill score values indicate an improvement of ECC-Q over the raw ensemble, while positive values indicate the contrary. Scenario D exhibits the smallest level of improvement, which is to expected, since scenario D reflects a situation where raw ensemble forecasts stem from the same distribution as the observations.

375 ~~CRPS skill score of raw ensemble with ECC-Q as reference model for the four different scenarios considered in Setting 3. The boxplots summarize results over 100 repetitions of each individual experiment.~~

## 4.2 Multivariate performance

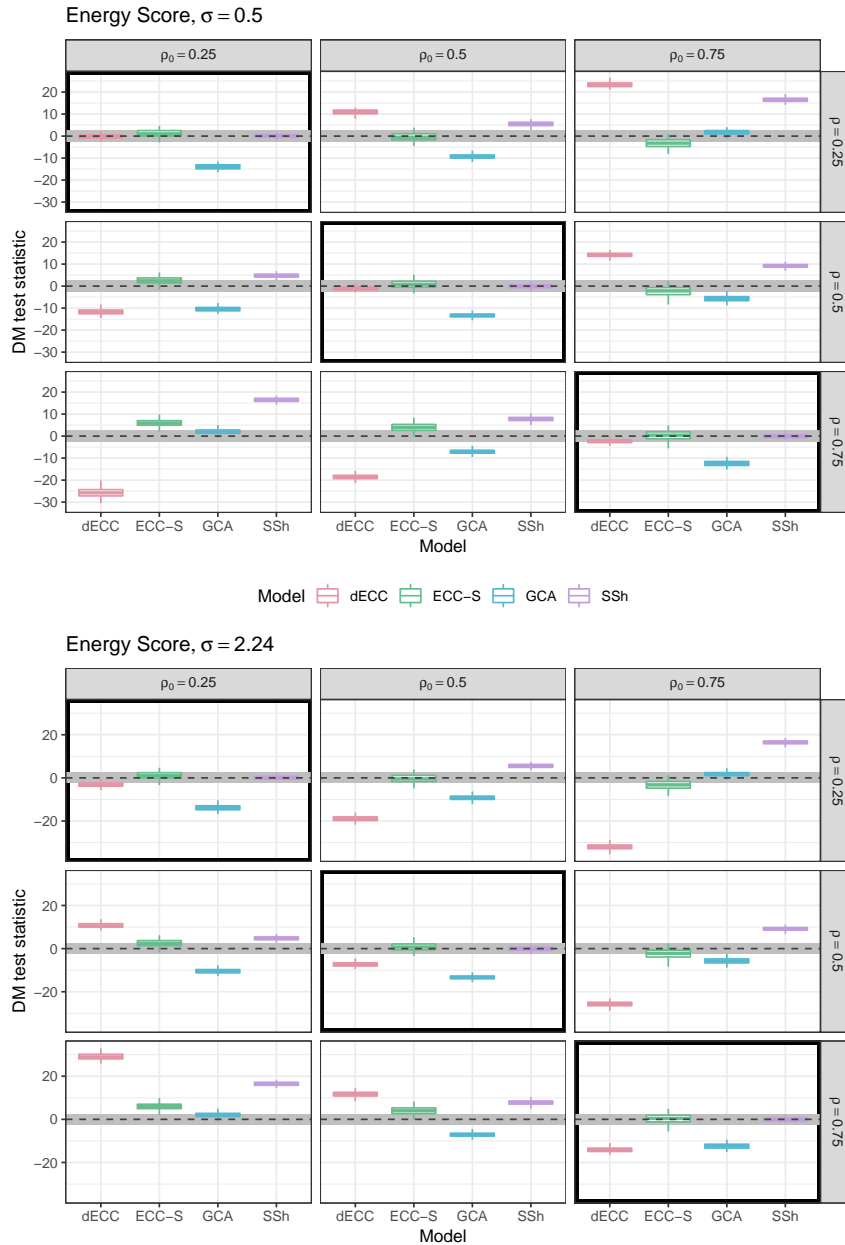
We now compare the multivariate performance of the different post-processing approaches presented in Section 2.2. Multivariate forecasts obtained by only applying the univariate post-processing methods without accounting for dependencies (denoted  
380 by EMOS-Q) as well as the raw ensemble predictions (ENS) are usually significantly worse and will be omitted in most comparisons below unless indicated otherwise. Additional figures with results for all parameter combinations in all settings are provided in the Supplementary Material.

### 4.2.1 Setting 1: Multivariate Gaussian distribution

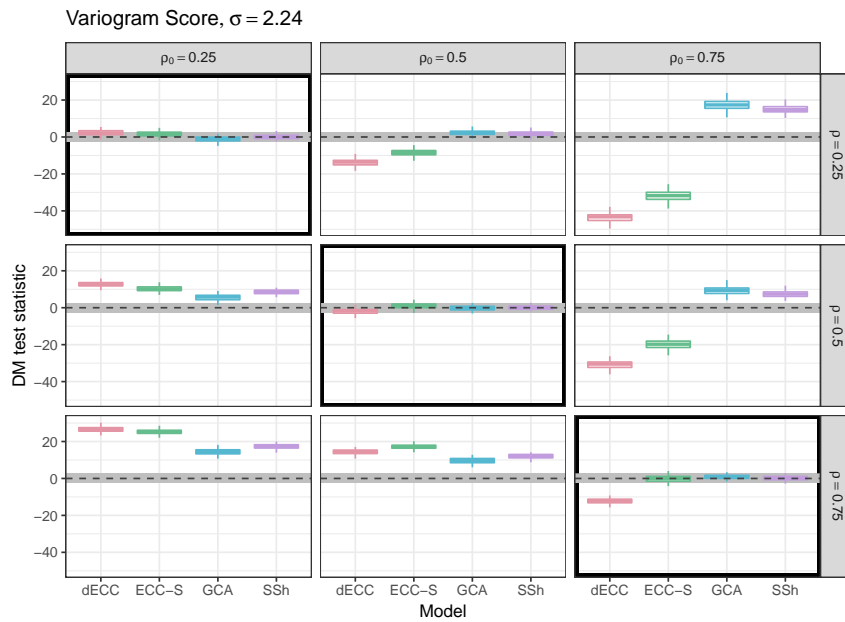
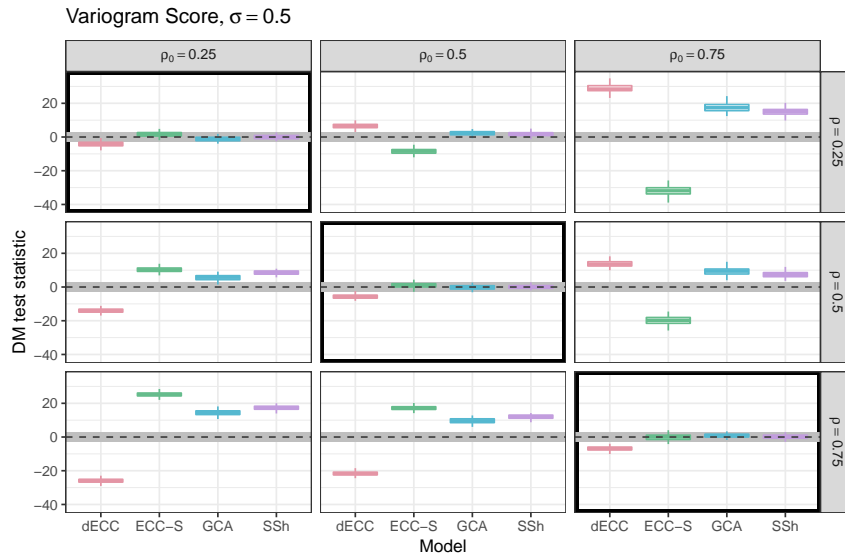
The tuning parameter  $\epsilon$  governing the bias in the mean vector of the ensemble forecasts only has very limited effects on the  
385 relative performance of the multivariate post-processing methods. To retain focus we restrict our attention to  $\epsilon = 1$ . Figure 2 shows results in terms of the ES for two different choices of  $\sigma$ , using multivariate forecasts of ECC-Q as reference method. For visual clarity, we omit parameter combinations where either  $\rho \in \{0.1, 0.9\}$  or  $\rho_0 \in \{0.1, 0.9\}$ . Corresponding results are available in the Supplementary Material. Note that the relative forecast performance of all approaches except for dECC generally does not depend on  $\sigma$ , ~~we~~. We thus proceed to discuss the remaining approaches first, and dECC last.

---

<sup>4</sup>~~In fact, ECC-Q does not change the marginal distributions, the univariate forecasts are thus identical to solely applying univariate post-processing methods in the margins separately, without accounting for dependencies. We will later refer to this as EMOS-Q.~~



**Figure 2.** Summaries of DM test statistic values based on the ES for Setting 1 with  $\epsilon = 1$ , and  $\sigma = 0.5$  (top), and  $\sigma = \sqrt{5}$  (bottom). ECC-Q forecasts are used as reference model such that positive values of the test statistic indicate improvements over ECC-Q and negative ~~value~~ indicated values indicate deterioration of forecast skill. Boxplots summarize results of the 100 Monte Carlo repetitions of each individual experiment. The horizontal gray stripe indicates the acceptance region of the two-sided DM test under the null hypothesis of equal predictive performance at a level of 0.05. Simulation parameter choices where the correlation structure of the raw ensemble is correctly specified ( $\rho = \rho_0$ ) are surrounded by black boxes.



**Figure 3.** As Figure 2, but summarizing results in terms of the VS.

390 If the correlation structure of the unprocessed ensemble forecasts is correctly specified (i.e.,  $\rho = \rho_0$ ), no significant differences can be detected between ECC-Q, ECC-S and SSh. In contrast, GCA (and dECC for larger values of  $\sigma$ ) perform substantially worse. The worse performance of GCA might be due to the larger forecast errors in the univariate margins, see Section 4.1.

In the cases with misspecifications in the correlation structure (i.e.,  $\rho \neq \rho_0$ ), larger differences can be detected among all 395 methods. Notably, SSh never performs substantially worse than ECC-Q and is always among the best performing approaches. This is not surprising as the only drawback of SSh in the present context and under the chosen implementation details is the underlying assumption of time-invariance of the correlation structure, which will be revisited in Setting 3. The larger the absolute difference between  $\rho$  and  $\rho_0$ , the greater the improvement of SSh relative to ECC-Q~~as~~. This is due to the fact that it becomes more and more beneficial to learn the dependence template from past observations rather than the raw ensemble, 400 the less information the ensemble provides about the true dependence structure. GCA also tends to outperform ECC-Q if the differences between  $\rho$  and  $\rho_0$  are large, however, GCA always performs worse than SSh and shows significantly worse performance than ECC-Q if the misspecifications in the ensemble are not too large (i.e., if  $\rho$  and  $\rho_0$  are equal or similar).

The relative performance of ECC-S depends on the ordering of  $\rho$  and  $\rho_0$ . If  $\rho > \rho_0$ , ECC-S significantly outperforms ECC-Q, however, if  $\rho < \rho_0$  significant ES differences in favor of ECC-Q can be detected. For dECC, the performance further depends 405 on the misspecification of the variance structure in the marginal distributions. If  $\rho > \rho_0$ , the DM test statistic values move from positive (improvement over ECC-Q) to negative (deterioration compared to ECC-Q) values for increasing  $\sigma$ . By contrast, if  $\rho < \rho_0$  the values of the test statistic instead change from negative to positive for increasing  $\sigma$ . The differences are mostly statistically significant, and indicate the largest relative improvements among all methods in cases of the largest possible differences between  $\rho$  and  $\rho_0$ . However, note that for some of those parameter combinations with small  $\rho$  and large  $\rho_0$ , even 410 EMOS-Q can outperform ECC-Q and ECC-S. In these situations, the raw ensemble forecasts contain very little information about the dependence structure and the ES can be improved by assuming independence instead of learning the dependence template from the ensemble.

Results in terms of the VS are shown in Figure 3. Most of the conclusions from the results in terms of the ES extend directly to comparisons based on the VS. SSh consistently remains among the best performing methods and provides significant im- 415 provements over ECC-Q unless  $\rho = \rho_0$ , however, alternative approaches here outperform SSh more often. Notably, the relative performance of GCA is consistently better in terms of the VS than in terms of the ES. For example, the differences between GCA and SSh appear to generally be negligible and GCA does not perform worse than ECC-Q for any of the simulation parameter combinations. These differences between the results for GCA in terms of ES and VS can potentially may be explained by the greater sensitivity of the VS to misspecifications in the correlation structure, whereas the ES shows a stronger dependence 420 on the mean vector.

For ECC-S and dECC, the general dependence on values of  $\rho$ ,  $\rho_0$  and  $\sigma$  (for dECC) is similar to the results for the ES, but the magnitude of both positive as well as negative differences to all other methods is increased. For example, it is now possible to find parameter combinations where either ECC-S or dECC (or both) substantially outperform both GCA and SSh.



## 4.2.2 Setting 2: Multivariate truncated Gaussian distributions

425 Summaries of DM test statistic values based on the ES (top) and the VS (bottom) for Setting 2 with  $c = 3$ ,  $\mu_0 = 2$ , and  $\sigma = 1$ .  
ECC-Q forecasts are used as reference model such that positive values of the test statistic indicate improvements over ECC-Q  
and negative value indicated deterioration of forecast skill. Boxplots summarize results of the 100 Monte Carlo repetitions of  
each individual experiment. The horizontal gray stripe indicates the acceptance region of the two-sided DM test under the null  
hypothesis of equal predictive performance at a level of 0.05. Simulation parameter choices where the correlation structure of  
430 the raw ensemble is correctly specified ( $\rho = \rho_0$ ) are surrounded by black boxes.

Figure ?? summarizes results for Setting 2 in terms of the ES (top) and the VS (bottom). In the interest of brevity, we only  
show results for  $\sigma = 1$ .

### The role of ensemble size $m$

To assess the effect of the ensemble size  $m$  on the results additional simulations have been performed with the simulation  
435 parameters from Figure 2, but ensemble sizes between 5 and  $\rho, \rho_0 \in \{0.25, 0.5, 0.75\}$ , but discuss the effect of misspecifications  
of the variance below. Corresponding plots 100. Corresponding figures are provided in the Supplemental Material.

Overall, SSH consistently provides significant improvements over ECC-Q except for settings in which the correlation  
structure of the ensemble is correctly specified ( $\rho = \rho_0$ ) where no significant differences can be detected. Overall, the relative  
440 ranking of the different methods is only very rarely effected by changes in the ensemble size. The relative differences in favor  
of SSH are increased for larger absolute differences of  $\rho$  and  $\rho_0$ . While these findings are in line with the results from Setting  
1 and hold for both ES and VS, the relative performance of GCA is now somewhat different from before. In particular terms  
terms of the ES between ECC-Q and ECC-S, and between ECC-Q and GCA become increasingly negligible with increasing  
ensemble size. Further, SSH shows improved predictive performance for larger numbers of ensemble members for  $\rho_0 < \rho$  in  
445 case of the ES, GCA here performs worse for many parameter settings. In particular, GCA is significantly worse than ECC-Q  
if  $\rho = \rho_0$ , and does not offer any improvements over ECC-Q if  $\rho$  and  $\rho_0$  are not too different. By contrast, in terms and for  
 $\rho_0 > \rho$  in case of the VSGCA shows much better performance and outperforms ECC-Q in almost all cases. Even for cases  
where  $\rho = \rho_0$ , no significant differences. The relative performance of dECC is strongly effected by changes in  $m$  for large  
misspecifications in the correlation parameters. A positive effect of larger numbers of members relative to ECC-Q in terms of  
445 VS both scoring rules can be detected.

450 As before, the results for ECC-S and dECC strongly depend on the misspecification of the correlation structure. Further, the  
results now additionally vary with the variance misspecification parameter for  $\rho_0 > \rho$  when  $\sigma < 1$ , and for  $\rho_0 < \rho$  when  $\sigma > 1$ .  
In both cases, the corresponding effects are negative if the misspecification in  $\sigma$  also for ECC-S is reversed.

### The role of dimension $d$

Additional simulations were further performed with dimensions  $d$  between 2 and 50 and the simulation parameters from above.  
455 In the interest of brevity, we refer to the Supplemental Material for corresponding figures. In terms of the ES, the results for

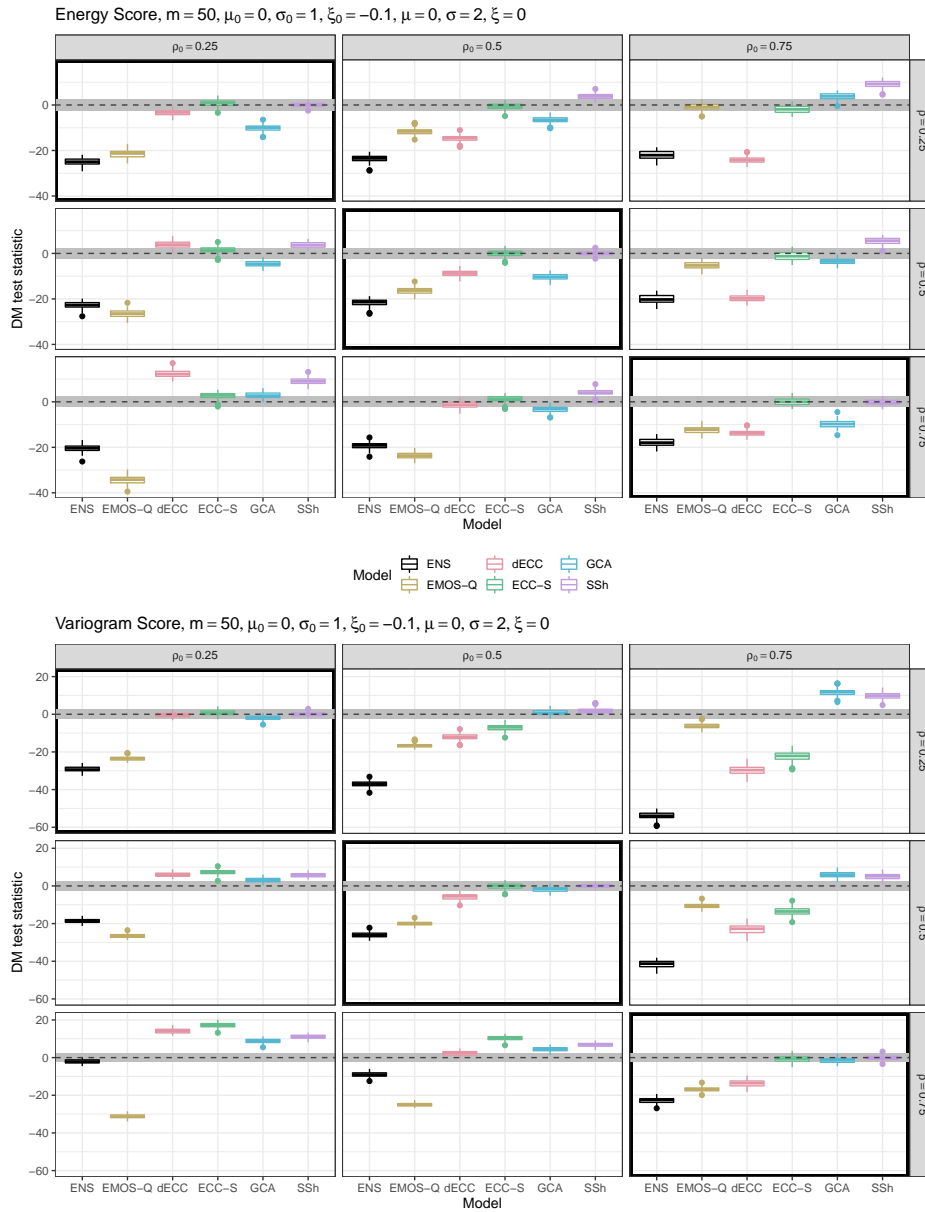
ECC-S ~~provides~~ are largely not effected by changes in dimension, whereas the relative performance of ECC-S improves with increasing  $d$  and minor improvements over ECC-Q for  $\rho > \rho_0$  and similar forecast quality for  $\rho \leq \rho_0$ . However, the forecast quality deteriorates for larger values of  $\sigma$  and the ECC-S forecasts are significantly worse than those of ECC-Q, even for cases where  $\rho = \rho_0$  can be detected even for correctly specified correlation parameters for high dimensions. For GCA, a marked deterioration of relative skill can be observed in terms of the ES, which can likely be attributed to sampling effects discussed above. In terms of the VS, ECC-S also provides significant improvements over ECC-Q for  $\rho > \rho_0$  and even provides the best forecasts among all models, but shows significantly worse performance for  $\rho < \rho_0$  across all values of  $\sigma$ .

The results for dECC are similar for ES and VS. For both scoring rules and  $\rho < \rho_0$ , the values of the DM test statistics move from positive (improvement over ECC-Q) to negative (deterioration compared to ECC-Q) values with increasing  $\sigma$ . For  $\rho > \rho_0$ , they instead change from negative to positive values. While improvements over ECC-Q can be observed for cases with a correctly specified correlation structure of the ensemble ( $\rho = \rho_0$ ) and  $\sigma = 1$ , dECC performs worse than ECC-Q for those cases if  $\sigma \neq 1$ . GCA partly shows the best relative performance among all methods for dimensions between 10 and 20, and performs worse in higher dimensions. The relative differences in predictive performance in favor of SSh are more pronounced in larger dimensions, in particular in cases with large misspecification of the correlation parameters. Changes of the relative performance of dECC in terms of both scoring rules for increasing numbers of dimensions are similar to those observed for increasing numbers of ensemble members.

#### 4.2.2 Setting 32: Multivariate censored GEV distributions

The four considered scenarios in Table 2 constitute different types of deviation of the ensemble from the observation properties. Results for Scenario B are given below in Figure 4, while the corresponding figures for Scenarios A, C, and D can be found in Section 2.1 of the Supplement. As the  $GEV_0$  distribution yields extreme outliers much more frequently than the Gaussian distribution in Setting 1, all figures (here and in the Supplemental Material) show only those values that are within  $1.5 \times$  interquartile range, so that the overall comparison of the boxplots does not suffer from single extreme outliers.

- In Scenario B the location is correctly specified, but scale and shape are misspecified such that ensemble forecasts have both larger scale and shape, resulting in a heavier right tail and slightly higher point masses at zero. This scenario is taken as a reference among the four considered ones and shown in Figure 4. Additional figures with results for the remaining scenarios are provided in the Supplemental Material. Multivariate post-processing improves considerably upon the raw ensemble. ECC-Q is outperformed by SSh and GCA only when the absolute difference between  $\rho_y$  and  $\rho_x$ ,  $\rho_0$ , and  $\rho$  becomes larger. As before, this is likely caused by the use of past observations to determine the dependence template by GCA and SSh which proves beneficial in comparison to ECC-Q in cases of a highly incorrect correlation structure in the ensemble. For correctly specified correlations (panels on the main diagonal in Figure 4), the relative performance of the methods does not depend on the actual value of correlation.
- In Scenario A the observation location parameter is shifted from 0 to a positive value for the ensemble, the observation scale is larger, and shape smaller than in the ensemble. Therefore, the ensemble forecasts come from a distribution



**Figure 4.** Summaries of DM test statistic values based on the ES (top) and the VS (bottom) for [setting Setting 2, scenario B](#) from Table 2, based on  $m = 50$  ensemble members. ECC-Q forecasts are used as reference model such that positive values of the test statistic indicate improvements over ECC-Q and negative value indicated values indicate deterioration of forecast skill. Boxplots summarize results of the 100 Monte Carlo repetitions of each individual experiment. The horizontal gray stripe indicates the acceptance region of the two-sided DM test under the null hypothesis of equal predictive performance at a level of 0.05.

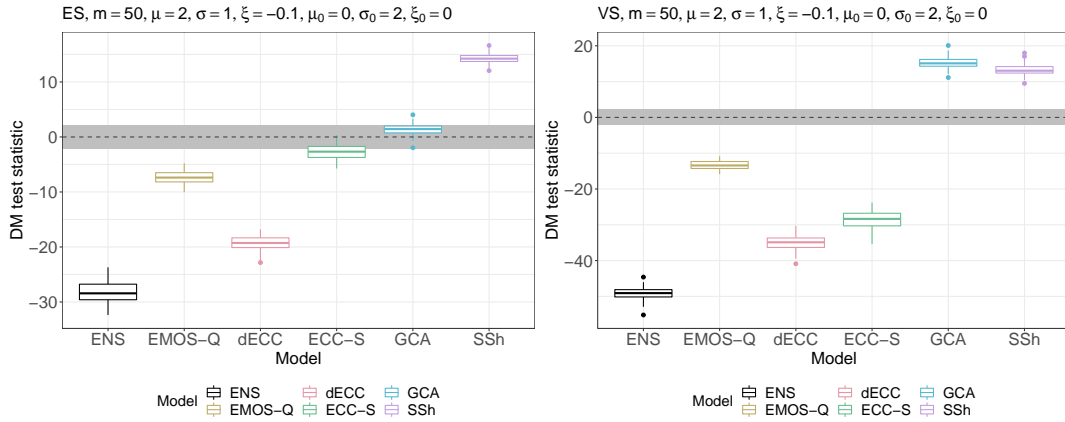
with smaller spread than the observations, which is also centered away from zero and has lower point mass at 0. In  
490 comparison to Scenario B there are more outliers, especially for ECC-S. In case of correctly specified correlations, the  
performance of the methods also does not depend on the actual value of correlation as in Scenario B. Notably, EMOS-  
Q here performs mostly similar to the ensemble, while in the other 3 scenarios it typically performs worse than the  
ensemble if  $\rho_x > \rho_y \rho_0 > \rho_0$ .

– In Scenario C the observation location is larger, the scale smaller, and the shape larger than in the ensemble distribution.  
495 This results in an observation distribution with a much more heavy right tail and a much larger point mass at 0 compared  
to the ensemble distribution. Here, post-processing models frequently offers no or only slight improvements over the  
raw ensemble. While ECC-Q does not always outperform the raw ensemble forecasts, SSh still shows improved forecast  
performance. As in the other scenarios, in case of correctly specified correlations, the performance of the methods does  
not depend on the actual value of correlation.

500 – In Scenario D all univariate distribution parameters are correctly specified. Therefore, the main differences in perfor-  
mance are imposed by the different misspecifications of the correlation structure. The main difference compared to the  
other scenarios is given by the markedly worse effects of not accounting for multivariate dependencies during post-  
processing (EMOS-Q).

In general, the methods perform differently across the four scenarios, but for most situations multivariate post-processing  
505 improves upon univariate post-processing without accounting for dependencies. Furthermore, SSh reveals a good performance  
in all four scenarios when  $\rho_y \rho_0$  differs considerably from  $\rho_x \rho_0$ . The performance of SSh has a tendency to improve further  
when the observation correlation is larger than the ensemble correlation. Within each of the four scenarios, the performance  
of the methods is nearly identical in cases where the correlation is correctly specified. In other words, as long as the ensemble  
forecasts correctly represent the correlation of the observations, the actual value of the correlation does not have an impact on  
510 the performance of a multivariate post-processing method. Above described observations can be found both in terms of the ES  
~~as well as~~ and the VS.

In addition to the scenarios from Table 2, further scenario variations were considered for  $\rho_y = 0.75$  and  $\rho_x = 0.25$   $\rho_0 = 0.75$   
and  $\rho = 0.25$ , that is for the case where ensemble correlation is too low compared to the observations. Figure 5 shows the  
situation where the observation location parameter is larger, the scale smaller, but the shape also smaller than in the ensemble  
515 forecasts. This contrasts the situation in Scenario C. While in C the observations were heavier tailed with higher point mass at 0,  
here it is the other way around (the ensemble distribution is heavier tailed with higher point mass at 0). In accordance ~~to~~  
with Scenarios A, B, C (where there are parameter misspecifications in the ensemble compared to the observations), EMOS-Q  
performs better than the raw ensemble and also better than dECC (as in B and C), while SSh and GCA perform best. However,  
by in contrast to results in terms of the ES, GCA exhibits an even better performance compared to the other models in terms of  
520 the VS. This indicates that the VS ~~might be~~ is better able to account for the correctly specified (or by post-processing improved)  
correlation structure than the ES.



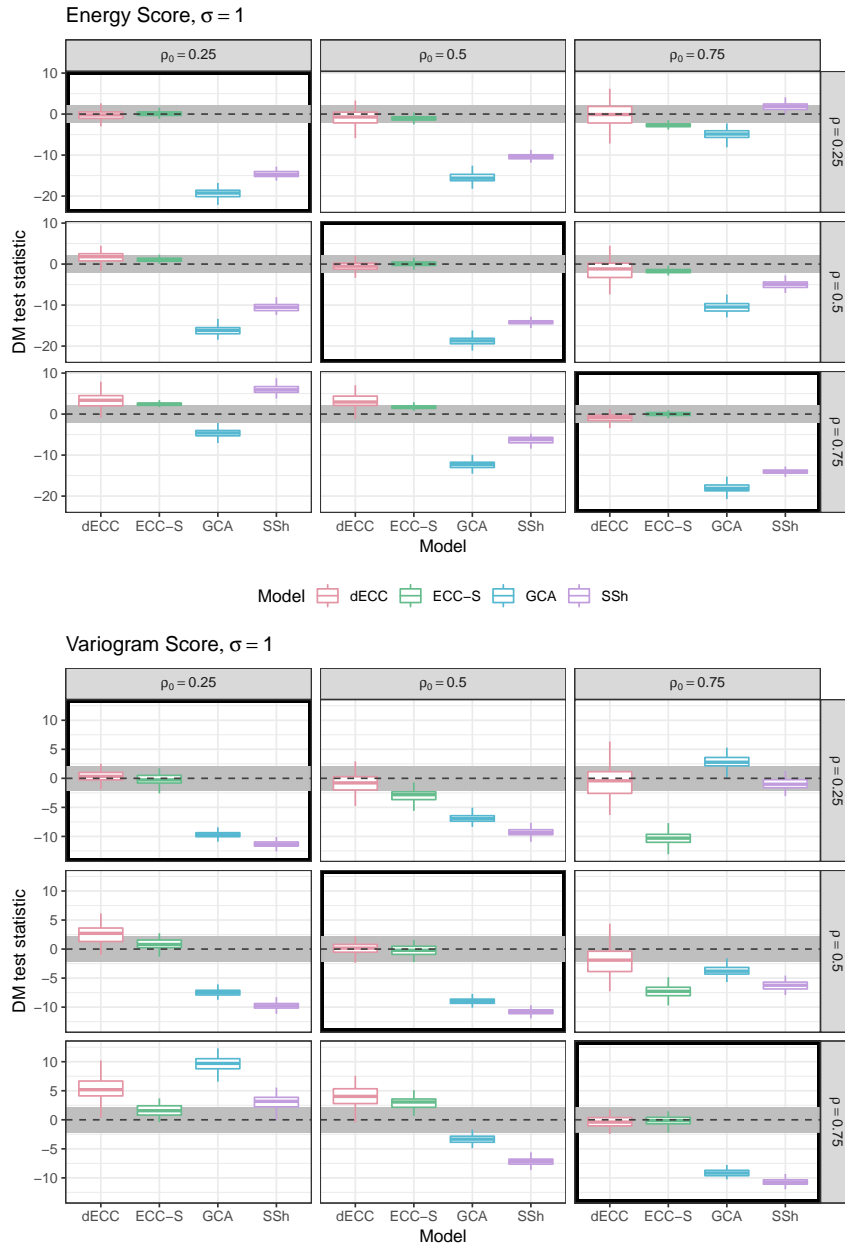
**Figure 5.** As Figure 4, but based on ES and VS for  $(\mu_y, \xi_y, \sigma_y) = (2.0, -0.1, 1.0)$   $(\mu_0, \xi_0, \sigma_0) = (2.0, -0.1, 1.0)$  and  $(\mu_x, \xi_x, \sigma_x) = (0.0, 0.0, 2.0)$   $(\mu, \xi, \sigma) = (0.0, 0.0, 2.0)$ , where  $\rho_y = 0.75$   $\rho_0 = 0.75$  and  $\rho_x = 0.25$   $\rho = 0.25$ , and ensemble size  $m = 50$

### The role of ensemble size $m$

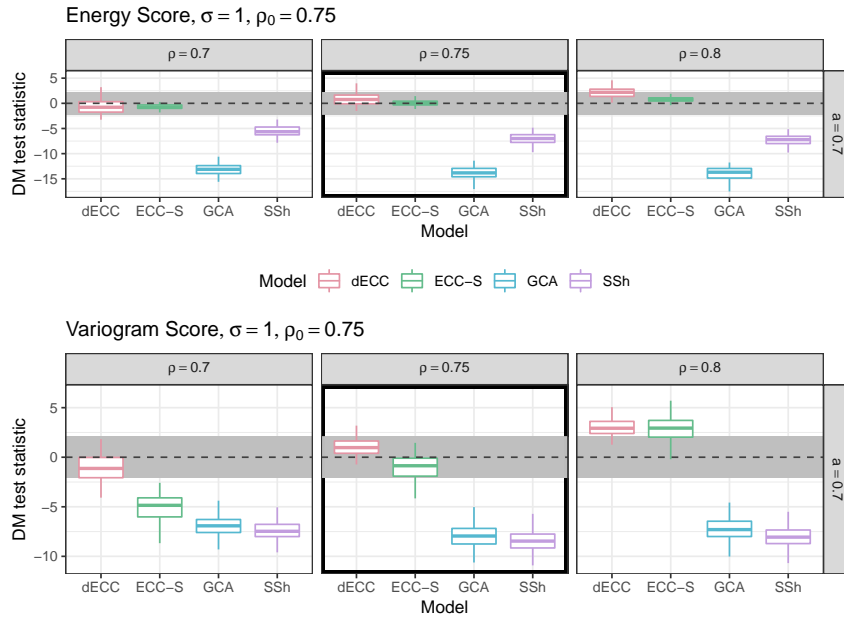
To assess the effect of the ensemble size  $m$  additional simulations have been performed for each of the four scenarios in Table 2 with ensemble sizes between 5 and 100. Corresponding comparative figures comparing ensemble sizes  $m = 5, 20, 50, 100$  for the Scenarios A, B, C, D are provided in Section 2.2 of the Supplemental Material. Overall, the size of the ensemble only has a minor effect on the relative performance of the multivariate methods apart from GCA, which strongly benefits from an increasing number of members across all four scenarios, specifically with regard to ES. This improvement is likely due to the sampling issues discussed above and is less pronounced in terms of the VS. As in Setting 1 the relative differences between ECC-Q and ECC-S in terms of the ES become increasingly negligible with increasing ensemble size in all considered scenarios (especially for  $\rho_0 = \rho$ ). This phenomenon is also less pronounced for the VS. On the contrary to the methods using the dependence information, the performance of EMOS-Q (not accounting for dependence) compared to ECC-Q becomes increasingly worse for increasing number of members when measured by ES. For VS, the influence of the number of members on EMOS-Q is only small. Interestingly, the difference in performance of the raw ensemble for an increasing number of members is negligible in case the misspecification is only minor and ES is considered. In case there is no misspecification (Scenario D), the raw ensemble can slightly benefit from an increasing number of members. Similar to the effect for ECC-Q, when measuring performance with VS, the effect on the raw ensemble is negligible. Further, it can be observed that the difference of the results between varying numbers of members is smallest for  $\rho_0 = \rho$ .

### 4.2.3 Setting 43: Multivariate Gaussian distribution with changes over iterations

Figure ??-6 shows results in terms of ES and VS for Setting 3A. We again only show results for  $\rho, \rho_0 \in \{0.25, 0.5, 0.75\}$  and refer to the Supplementary Material for further results. The most notable differences compared to Setting 1 are that the different



**Figure 6.** Summaries of DM test statistic values based on the ES (top) and the VS (bottom) for Setting 4-3A with  $\epsilon = 1$  and  $\sigma = 1$ . ECC-Q forecasts are used as reference model such that positive values of the test statistic indicate improvements over ECC-Q and negative ~~value~~ indicated values indicate deterioration of forecast skill. Boxplots summarize results of the 100 Monte Carlo repetitions of each individual experiment. The horizontal gray stripe indicates the acceptance region of the two-sided DM test under the null hypothesis of equal predictive performance at a level of 0.05. Simulation parameter choices where the correlation structure of the raw ensemble is correctly specified ( $\rho = \rho_0$ ) are surrounded by black boxes.



**Figure 7.** Summaries of DM test statistic values based on the ES (top) and the VS (bottom) for Setting 3B for  $\sigma = 1$  and high values of  $\rho, \rho_0$ . ECC-Q forecasts are used as reference model such that positive values of the test statistic indicate improvements over ECC-Q and negative values indicate deterioration of forecast skill. Boxplots summarize results of the 100 Monte Carlo repetitions of each individual experiment. The horizontal gray stripe indicates the acceptance region of the two-sided DM test under the null hypothesis of equal predictive performance at a level of 0.05. Simulation parameter choices where the correlation structure of the raw ensemble is correctly specified ( $\rho = \rho_0$ ) are surrounded by black boxes.

ECC variants here significantly outperform GCA and SSh not only for ensemble forecasts with correctly specified correlation structure, but also for small deviations of  $\rho$  from  $\rho_0$ . Significant ES differences in favor of SSh are only obtained for large absolute differences of  $\rho$  and  $\rho_0$ . Similar observations hold for GCA which, however, generally exhibits worse performance compared to SSh. The ES differences among the ECC variants are only minor and usually not statistically significant.

545 Similar conclusions apply for the VS, however, GCA generally performs better than SSh, and ECC-S provides significantly worse forecasts compared to the other ECC variants for  $\rho < \rho_0$ .

550 Results for Setting 3B are shown in Figure 7. Note that the columns here show different values of  $\rho$  and the row refers to a specific value of  $a$ . Similar to Setting 3A, we observe that in terms of the ES, dECC and ECC-S do not show significant differences in performance compared to ECC-Q, whereas GCA and SSh here perform worse for all parameter combinations. In terms of the VS, also GCA now performs worse than ECC-Q for all correlation parameters, whereas significantly negative and positive differences of ECC-S and dECC compared to ECC-Q can be detected for  $\rho < \rho_0$  and  $\rho > \rho_0$ , respectively. Additional results for varying values of  $\sigma$  and  $a$ , and sets of low or medium correlations are provided in the Supplemental Material. The results generally do not depend on the choice of  $\sigma$ . Results for low and medium correlation parameter values are characterized

555 by less substantial differences between the methods. In particular, it is only rarely possible to detect significant differences  
when comparing ECC-Q and SSh, and GCA only performs significantly worse in terms of the ES. Further, there exist more  
parameter combinations with improvements by ECC-S and dECC. However, note that due to the setup of Setting 3B, the  
variations over time in both observations and ensemble predictions will be much smaller than for high correlation parameter  
values. Within a fixed set of correlation parameters, the relative differences between the methods become more pronounced  
with increasing values of  $a$ .

560 Note that the main focus here in both variants of Setting 3 was to demonstrate that in (potentially more realistic) settings  
with changes over iterationstime, naive implementations of the Schaake shuffle can perform worse than ECC variants. However,  
similarity-based implementations of the Schaake shuffle (Schefzik, 2016; Scheuerer et al., 2017) are available and may be able  
to alleviate this issue.

## 5 Discussion and conclusion

565 State of the art methods for multivariate ensemble post-processing were compared in four simulation settings which aimed to  
mimic different situations and challenges occurring in practical applications. Across all settings, the Schaake shuffle consti-  
tutes a powerful benchmark method that proves difficult to outperform, except for naive implementations in the presence of  
structural change (for example, time-varying correlation structures considered in Setting 3). By contrast to SSh, the Gaussian  
copula approach typically only provides improvements over variants of ensemble copula coupling if the parametric assump-  
570 tion of a Gaussian copula is satisfied or if forecast performance is evaluated with the variogram score. Results in terms of the  
CRPS further highlight an additional potential disadvantage in that the univariate forecast errors are larger compared to the  
competitors.

Not surprisingly, variants of ensemble copula coupling typically perform the better the more informative the ensemble  
forecasts are about the true multivariate dependence structure. A particular advantage compared to standard implementations of  
575 SSh and GCA illustrated in Setting 4-3 may be given by the ability to account for flow-dependent differences in the multivariate  
dependence structure if those are (at least approximately) present in the ensemble predictions, but not in a randomly selected  
subset of past observations.

There is no consistently best method across all simulation settings and potential misspecifications among the different ECC  
variants investigated here (ECC-Q, ECC-S and dECC). ECC-Q provides a reasonable benchmark model and will rarely yield  
580 the worst forecasts among all ECC variants. Significant improvements over ECC-Q may be obtained by ECC-S and dECC ;  
however in specific situations, including specific combinations of ensemble size and dimension. For example, dECC sometimes  
works well for underdispersive ensembles where the correlation is too low, whereas ECC-S may work better if the ensemble  
is underdispersive and the correlation is too strong. However, the results will strongly depend on the exact misspecification of  
the variance-covariance structure of the ensemble as well as the performance measure chosen for multivariate evaluation.

585 In light of the presented results it seems to be generally advisable to first test the Schaake shuffle along with ECC-Q. If  
structural assumptions on specific misspecifications of the ensemble predictions seem appropriate, extensions by other variants



of ECC or GCA might provide improvements. However, it should be noted that the results for real-world ensemble prediction systems may be influenced by many additional factors, and may differ when considering station-based or grid-based post-processing methods. The computational costs of all presented methods are not only negligible in comparison to the generation of the raw ensemble forecasts, but also compared to the univariate post-processing as no numerical optimization is required. It may thus be generally advisable to compare multiple multivariate post-processing methods for the specific dataset and application at hand.

The simulation settings considered here provide several avenues for further generalization and analysis. For example, a ~~detailed investigation of the effect of the dimension  $d$  and a~~ comparison of forecast quality in terms of multivariate calibration (Thorarinsdottir et al., 2016; Wilks, 2017) is left for future work. Further, the autoregressive structure of the correlations across dimensions may be extended towards more complex correlation functions, see, e.g., Thorarinsdottir et al. (2016, Section 4.2). While we only considered multivariate methods based on a two step procedure combining univariate post-processing and dependence modeling via copulas, an extension of the comparison to parametric approaches along the lines of Feldmann et al. (2015) and Baran and Möller (2015) present another starting point for future work. Note that within the specific choices for Setting 1, the spatial EMOS approach of Feldmann et al. (2015) can be seen as a special case of GCA.

We have limited our investigation to simulation studies only as those settings allow to readily assess the effects of different types of misspecifications of the various multivariate properties of ensemble forecasts and observations, and may thus help to guide implementations of multivariate post-processing. ~~An~~ Further, they are able to provide a more complete picture of the effects of different types of misspecifications on the performance of the different methods than those that may be observed in practical applications. Nonetheless, an important aspect for future work is to complement the comparison of multivariate post-processing methods by studies based on real-world datasets of ensemble forecasts and observations, extending existing comparisons of subsets of the methods considered here (e.g., Schefzik et al., 2013; Wilks, 2015). However, the variety of application scenarios, methods and implementation choices likely requires large-scale efforts, ideally based on standardized benchmark datasets. A possible intermediate step might be given by the use of simulated datasets obtained via stochastic weather generators (see, e.g., Wilks and Wilby, 1999) which may provide arbitrarily large datasets with possibly more realistic properties than the simple settings considered here.

A different perspective on the results presented here concerns the evaluation of multivariate probabilistic forecasts. In recent work Ziel and Berk (2019) argue that the use of Diebold-Mariano tests is of crucial importance for appropriately assessing the discrimination ability of multivariate proper scoring rules and find that the ES might not have as bad discrimination ability as indicated by earlier research. The simulation settings and comparisons of multivariate post-processing methods considered here may be seen as additional simulation studies for assessing the discrimination ability of multivariate proper scoring rules. In particular, the results in Section 4 are in line with the findings of Ziel and Berk (2019) in that the ES does not exhibit inferior discrimination ability compared to the VS. Nonetheless, the ranking of the different multivariate post-processing methods strongly depends on the proper scoring rule used for evaluation, and further research on multivariate verification is required to address open questions, improve mathematical understanding and guide model comparisons in applied work.

*Code availability.* R code with implementations of all simulation settings as well as code to reproduce the results presented here and in the Supplemental Material is available from [https://github.com/slerch/multiv\\_pp](https://github.com/slerch/multiv_pp).

## Appendix A: Details on the left-censored generalized extreme value (GEV) distribution

When the GEV distribution is left-censored at zero, its cumulative distribution function can be written as

$$F(y) = \begin{cases} e^{-t(y)}, & y \geq 0 \\ 0, & y < 0 \end{cases}, \quad \text{where } t(y) = \begin{cases} (1 + \xi \left(\frac{y-\mu}{\sigma}\right))^{-1/\xi}, & \xi \neq 0 \\ e^{-(y-\mu)/\sigma}, & \xi = 0 \end{cases}$$

for  $y \in \mathfrak{Y}$ , where  $\mathfrak{Y} = [\mu - \sigma/\xi, \infty)$  when  $\xi > 0$ ,  $\mathfrak{Y} = (-\infty, \infty)$  when  $\xi = 0$  and  $\mathfrak{Y} = (-\infty, \mu - \sigma/\xi]$  when  $\xi < 0$ . This describes a three-parameter distribution family, where  $\mu \in \mathbb{R}$ ,  $\sigma > 0$ , and  $\xi \in \mathbb{R}$  are location, scale, and shape of the non-censored GEV distribution, respectively.

### Expectation and variance

Let  $Y$  be a random variable distributed according to GEV and censored at zero to the left. From the law of total expectation

$$\mathbb{E}(g(Y)) = P(Y = 0)\mathbb{E}(g(Y)|Y = 0) + P(Y > 0)\mathbb{E}(g(Y)|Y > 0),$$

where the second term in the sum is given by

$$\mathbb{E}(g(Y)\mathbb{1}_{\{Y>0\}}) = \int_0^{\infty} g(y)f_Y(y)dy.$$

Here,  $\mathbb{1}$  denotes the indicator function,  $g$  is any function of  $Y$  such that  $g(Y)$  is a random variable, and  $f_Y$  is the probability density function (PDF) of the non-censored GEV. By noting that  $\mathbb{E}(Y|Y = 0) = \mathbb{E}(Y^2|Y = 0) = 0$ , expectation and variance of the left-censored GEV can be computed from the two integrals  $\int_0^{\infty} yf_Y(y)dy$  and  $\int_0^{\infty} y^2f_Y(y)dy$ , the former existing when  $\xi < 1$  and the latter existing when  $\xi < 0.5$ . Both integrals are not derived analytically here, but evaluated by numerical integration. In contrast to the non-censored GEV distribution, the variance of the left-censored version also depends on the parameter  $\mu$ , since different choices of  $\mu$  lead to different left-censored CDFs which are not merely distinguished by location. Therefore  $\mu$  is a location parameter for the non-censored GEV, but not for the left-censored version.

## Appendix B: Evaluating probabilistic forecasts

### B1 Proper scoring rules

The comparative evaluation of probabilistic forecasts is usually based on proper scoring rules. A proper scoring rule is a function

$$S : \mathcal{F} \times \Omega \rightarrow \mathbb{R},$$

640 which assigns a numerical score  $S(F, y)$  to a pair of a forecast distribution  $F \in \mathcal{F}$  and a realizing observation  $y \in \Omega$ . Here,  $\mathcal{F}$  denotes a class of probability distributions supported on  $\Omega$ . The forecast distribution  $F$  may come in the form of a predictive CDF, PDF, or a discrete sample as in the case of ensemble predictions. A scoring rule is called proper if

$$\mathbb{E}_G S(G, Y) \leq \mathbb{E}_G S(F, Y)$$

for all  $F, G \in \mathcal{F}$ , and strictly proper if equality holds only if  $F = G$ . See Gneiting and Raftery (2007) for a review of proper  
645 scoring rules from a statistical perspective.

The most popular example of a univariate (i.e.,  $\Omega \subset \mathbb{R}$ ) proper scoring rule in the environmental sciences is given by the continuous ranked probability score (CRPS),

$$\text{CRPS}(F, y) = \int_{\Omega} (F(z) - \mathbb{1}\{z \geq y\})^2 dz.$$

Over the past years a growing interest in multivariate proper scoring rules accompanies the proliferation of multivariate  
650 probabilistic forecasting methods in applications across disciplines. The definition of proper scoring rules from above straightforwardly extends towards multivariate settings (i.e.,  $\Omega \subset \mathbb{R}^d$ ). A variety of multivariate proper scoring rules has been proposed over the past years, usually focused on cases where multivariate probabilistic forecasts are given as samples from the forecast distributions.

To introduce multivariate scoring rules let  $\mathbf{y} = (y^{(1)}, \dots, y^{(d)}) \in \Omega \subset \mathbb{R}^d$ , and let  $F$  denote a forecast distribution on  $\mathbb{R}^d$   
655 given through by  $m$  discrete samples  $\mathbf{X}_1, \dots, \mathbf{X}_m$  from  $F$  with  $\mathbf{X}_i = (X_i^{(1)}, \dots, X_i^{(d)}) \in \mathbb{R}^d, i = 1, \dots, m$ . Important examples of multivariate proper scoring rules include the energy score (ES; Gneiting et al., 2008),

$$\text{ES}(F, y) = \frac{1}{m} \sum_{i=1}^m \|\mathbf{X}_i - \mathbf{y}\| - \frac{1}{2m^2} \sum_{i=1}^m \sum_{j=1}^m \|\mathbf{X}_i - \mathbf{X}_j\|,$$

where  $\|\cdot\|$  is the Euclidean norm on  $\mathbb{R}^d$ , and the variogram score of order  $p$  ( $\text{VS}^p$ ; Scheuerer and Hamill, 2015),

$$\text{VS}^p(F, y) = \sum_{i=1}^d \sum_{j=1}^d w_{i,j} \left( \left| y^{(i)} - y^{(j)} \right|^p - \frac{1}{m} \sum_{k=1}^m \left| X_k^{(i)} - X_k^{(j)} \right|^p \right)^2.$$

660 Here,  $w_{i,j}$  is a non-negative weight that allows to emphasize or down-weight pairs of component combinations, and  $p$  is the order of the variogram score. Following suggestions of Scheuerer and Hamill (2015), we considered  $p = 0.5$  and  $p = 1$ . As none of the simulations settings indicated any substantial differences, we set  $p = 1$  throughout and denote  $\text{VS}^1(F, y)$  by  $\text{VS}(F, y)$ . Since the generic multivariate structure of the simulation settings does not impose any meaningful structure in pairs of components we focus on the unweighted versions of the variogram score. Several weighting schemes have been tested, but  
665 did not lead to any substantially different conclusions.

We utilize implementations provided in the R package `scoringRules` (Jordan et al., 2019) to compute univariate and multivariate scoring rules for forecast evaluation and post-processing model estimation.

## B2 Diebold-Mariano tests

670 Statistical tests of equal predictive performance are frequently used to assess the statistical significance of observed score differences between models. We focus on Diebold-Mariano (DM; Diebold and Mariano, 1995) tests which are widely used in the econometric literature due to their ability to account for temporal dependencies. For applications in the context of post-processing, see, e.g., Baran and Lerch (2016).

For a (univariate or multivariate) proper scoring rule  $S$  and sets of two competing probabilistic forecasts  $F_i$  and  $G_i$ ,  $i = 1, \dots, n_{\text{test}}$  over a test set, the test statistic of the DM test is given by

$$675 \quad T_{n_{\text{test}}}^{\text{DM}} = \frac{\overline{S(F, y)} - \overline{S(G, y)}}{\hat{\sigma}}, \quad (\text{B1})$$

where  $\overline{S(F, y)} = \frac{1}{n_{\text{test}}} \sum_{i=1}^{n_{\text{test}}} S(F_i, y_i)$  and  $\overline{S(G, y)} = \frac{1}{n_{\text{test}}} \sum_{i=1}^{n_{\text{test}}} S(G_i, y_i)$  denote the mean score values of  $F$  and  $G$  over the test set of size  $n_{\text{test}}$ , respectively. In (B1),  $\hat{\sigma}$  denotes an estimator of the asymptotic standard deviation of the sequence of score differences of  $F$  and  $G$ . Positive values of  $T_{n_{\text{test}}}^{\text{DM}}$  indicate a superior performance of  $G$ , whereas negative values **indicated** indicate a superior performance of  $F$ .

680 Under standard regularity assumptions and the null hypothesis of equal predictive performance,  $T_{n_{\text{test}}}^{\text{DM}}$  asymptotically follows a standard normal distribution which allows to assess the statistical significance of differences in predictive performance. We utilize implementations of DM tests provided in the R package `forecast` (Hyndman and Khandakar, 2008).

*Author contributions.* All authors jointly discussed and devised the design and setup of the simulation studies. A variant of Setting 1 was first investigated in a MSc thesis written by MG (Graeter, 2016), co-supervised by SL. SL wrote evaluation and plotting routines, implemented simulation settings 1 and 4, partially based on code and suggestions from MG and SH, and provided a simulation framework in which the variant of Setting 1 based on a multivariate truncated normal distribution (SB) and Setting 2 (AM and JG) were implemented. All authors jointly analyzed the results and edited the manuscript, coordinated by SL.

*Competing interests.* Sebastian Lerch and Stephan Hemri are editors of the special issue on “Advances in post-processing and blending of deterministic and ensemble forecasts”. The remaining authors declare that they have no conflict of interest.

690 *Acknowledgements.* The authors gratefully acknowledge support by the Deutsche Forschungsgemeinschaft (DFG) through project MO-3394/1-1 “Statistische Nachbearbeitung von Ensemble-Vorhersagen für verschiedene Wettervariablen”. Sebastian Lerch is further supported by DFG through SFB/TRR 165 “Waves to Weather”, and Sándor Baran by the National Research, Development and Innovation Office under Grant No. NN125679. The authors thank Tilmann Gneiting and Kira Feldmann for helpful discussions. [Constructive comments on an earlier version of the manuscript by Zied Ben Bouallègue and two anonymous referees are gratefully acknowledged.](#)

- Allen, S., Ferro, C. A. T., and Kwasniok, F.: Regime-dependent statistical post-processing of ensemble forecasts, *Quarterly Journal of the Royal Meteorological Society*, Early View, <https://doi.org/10.1002/qj.3638>, 2019.
- Baran, S. and Lerch, S.: Mixture EMOS model for calibrating ensemble forecasts of wind speed, *Environmetrics*, 27, 116–130, <https://doi.org/10.1002/env.2380>, 2016.
- 700 Baran, S. and Möller, A.: Joint probabilistic forecasting of wind speed and temperature using Bayesian model averaging, *Environmetrics*, 26, 120–132, <https://doi.org/10.1002/env.2316>, 2015.
- Baran, S. and Möller, A.: Bivariate ensemble model output statistics approach for joint forecasting of wind speed and temperature, *Meteorology and Atmospheric Physics*, 129, 99–112, <https://doi.org/10.1007/s00703-016-0467-8>, 2017.
- Ben Bouallègue, Z., Heppelmann, T., Theis, S. E., and Pinson, P.: Generation of Scenarios from Calibrated Ensemble Forecasts with a Dual-  
705 Ensemble Copula-Coupling Approach, *Monthly Weather Review*, 144, 4737–4750, <https://doi.org/10.1175/MWR-D-15-0403.1>, 2016.
- Berrocal, V. J., Raftery, A. E., Gneiting, T., et al.: Probabilistic quantitative precipitation field forecasting using a two-stage spatial model, *The Annals of Applied Statistics*, 2, 1170–1193, <https://doi.org/10.1214/08-AOAS203>, 2008.
- Bröcker, J.: Evaluating raw ensembles with the continuous ranked probability score, *Quarterly Journal of the Royal Meteorological Society*, 138, 1611–1617, <https://doi.org/10.1002/qj.1891>, 2012.
- 710 Cario, M. C. and Nelson, B. L.: Modeling and generating random vectors with arbitrary marginal distributions and correlation matrix, Tech. rep., Department of Industrial Engineering and Management Sciences, Northwestern University, 1997.
- Chaloulos, G. and Lygeros, J.: Effect of wind correlation on aircraft conflict probability, *Journal of Guidance, Control, and Dynamics*, 30, 1742–1752, <https://doi.org/10.2514/1.28858>, 2007.
- Chen, H.: Initialization for NORTA: Generation of random vectors with specified marginals and correlations, *INFORMS Journal on Computing*, 13, 312–331, <https://doi.org/10.1287/ijoc.13.4.312.9736>, 2001.
- 715 Clark, M., Gangopadhyay, S., Hay, L., Rajagopalan, B., and Wilby, R.: The Schaake shuffle: A method for reconstructing space–time variability in forecasted precipitation and temperature fields, *Journal of Hydrometeorology*, 5, 243–262, [https://doi.org/10.1175/1525-7541\(2004\)005<0243:tssamf>2.0.co;2](https://doi.org/10.1175/1525-7541(2004)005<0243:tssamf>2.0.co;2), 2004.
- Diebold, F. X. and Mariano, R. S.: Comparing predictive accuracy, *Journal of Business and Economic Statistics*, 13, 253–263,   
720 <https://doi.org/10.1198/073500102753410444>, 1995.
- Feldmann, K., Scheuerer, M., and Thorarindottir, T. L.: Spatial postprocessing of ensemble forecasts for temperature using nonhomogeneous Gaussian regression, *Monthly Weather Review*, 143, 955–971, <https://doi.org/10.1175/MWR-D-14-00210.1>, 2015.
- Feng, S., Nadarajah, S., and Hu, Q.: Modeling annual extreme precipitation in China using the generalized extreme value distribution, *Journal of the Meteorological Society of Japan. Ser. II*, 85, 599–613, <https://doi.org/10.2151/jmsj.85.599>, 2007.
- 725 Gneiting, T. and Raftery, A. E.: Strictly Proper Scoring Rules, Prediction, and Estimation, *Journal of the American Statistical Association*, 102, 359–378, <https://doi.org/10.1198/016214506000001437>, 2007.
- Gneiting, T., Raftery, A. E., Westveld III, A. H., and Goldman, T.: Calibrated probabilistic forecasting using ensemble model output statistics and minimum CRPS estimation, *Monthly Weather Review*, 133, 1098–1118, <https://doi.org/10.1175/MWR2904.1>, 2005.
- Gneiting, T., Stanberry, L. I., Grimit, E. P., Held, L., and Johnson, N. A.: Assessing Probabilistic Forecasts of Multivariate Quantities, with  
730 an Application to Ensemble Predictions of Surface Winds, *Test*, 17, 211–235, <https://doi.org/10.1007/s11749-008-0114-x>, 2008.
- Graeter, M.: Simulation study of dual ensemble copula coupling, Master’s thesis, Karlsruhe Institute of Technology, 2016.

- Hemri, S. and Klein, B.: Analog-Based Postprocessing of Navigation-Related Hydrological Ensemble Forecasts, *Water Resources Research*, 53, 9059–9077, <https://doi.org/10.1002/2017WR020684>, 2017.
- 735 Hu, Y., Schmeits, M. J., van Andel, J. S., Verkade, J. S., Xu, M., Solomatine, D. P., and Liang, Z.: A Stratified Sampling Approach for Improved Sampling from a Calibrated Ensemble Forecast Distribution, *Journal of Hydrometeorology*, 17, 2405–2417, <https://doi.org/10.1175/JHM-D-15-0205.1>, 2016.
- Hyndman, R. J. and Khandakar, Y.: Automatic time series forecasting: the forecast package for R, *Journal of Statistical Software*, 26, 1–22, <https://doi.org/10.18637/jss.v027.i03>, 2008.
- 740 Jenkinson, A. F.: The frequency distribution of the annual maximum (or minimum) values of meteorological elements, *Quarterly Journal of the Royal Meteorological Society*, 81, 158–171, <https://doi.org/10.1002/qj.49708134804>, 1955.
- Jordan, A., Krüger, F., and Lerch, S.: Evaluating Probabilistic Forecasts with scoringRules, *Journal of Statistical Software*, 90, 1–37, <https://doi.org/10.18637/jss.v090.i12>, 2019.
- Lang, M. N., Mayr, G. J., Stauffer, R., and Zeileis, A.: Bivariate Gaussian models for wind vectors in a distributional regression framework, *Advances in Statistical Climatology, Meteorology and Oceanography*, 5, 115–132, <https://doi.org/10.5194/ascmo-5-115-2019>, 2019.
- 745 Lang, M. N., Lerch, S., Mayr, G. J., Simon, T., Stauffer, R., and Zeileis, A.: Remember the past: A comparison of time-adaptive training schemes for non-homogeneous regression, *Nonlinear Processes in Geophysics*, 27, 23–34, <https://doi.org/10.5194/npg-2019-49>, 2020.
- Lerch, S. and Thorarindottir, T. L.: Comparison of non-homogeneous regression models for probabilistic wind speed forecasting, *Tellus A*, 65, 21–206, <https://doi.org/10.3402/tellusa.v65i0.21206>, 2013.
- 750 Möller, A., Lenkoski, A., and Thorarindottir, T. L.: Multivariate probabilistic forecasting using ensemble Bayesian model averaging and copulas, *Quarterly Journal of the Royal Meteorological Society*, 139, 982–991, <https://doi.org/10.1002/qj.2009>, 2013.
- Morrison, J. E. and Smith, J. A.: Stochastic modeling of flood peaks using the generalized extreme value distribution, *Water Resources Research*, 38, 41–1–41–12, <https://doi.org/10.1029/2001WR000502>, 2002.
- Nelsen, R. B.: *An Introduction to Copulas*, Springer, New York, 2nd edn., 2006.
- 755 Pinson, P. and Girard, R.: Evaluating the quality of scenarios of short-term wind power generation, *Applied Energy*, 96, 12–20, <https://doi.org/10.1016/j.apenergy.2011.11.004>, 2012.
- Pinson, P. and Messner, J. W.: Application of postprocessing for renewable energy, in: *Statistical Postprocessing of Ensemble Forecasts*, edited by Vannitsem, S., Wilks, D. S., and Messner, J. W., pp. 241–266, Elsevier, 2018.
- R Core Team: *R: A Language and Environment for Statistical Computing*, R Foundation for Statistical Computing, Vienna, Austria, <https://www.R-project.org/>, 2019.
- 760 Raftery, A. E., Gneiting, T., Balabdaoui, F., and Polakowski, M.: Using Bayesian model averaging to calibrate forecast ensembles, *Monthly Weather Review*, 133, 1155–1174, <https://doi.org/10.1175/MWR2906.1>, 2005.
- Rasp, S. and Lerch, S.: Neural Networks for Postprocessing Ensemble Weather Forecasts, *Monthly Weather Review*, 146, 3885–3900, <https://doi.org/10.1175/MWR-D-18-0187.1>, 2018.
- 765 Schefzik, R.: A similarity-based implementation of the Schaake shuffle, *Monthly Weather Review*, 144, 1909–1921, <https://doi.org/10.1175/MWR-D-15-0227.1>, 2016.
- Schefzik, R.: Ensemble calibration with preserved correlations: unifying and comparing ensemble copula coupling and member-by-member postprocessing, *Quarterly Journal of the Royal Meteorological Society*, 143, 999–1008, <https://doi.org/10.1002/qj.2984>, 2017.
- Schefzik, R. and Möller, A.: Ensemble postprocessing methods incorporating dependence structures, in: *Statistical Postprocessing of Ensemble Forecasts*, edited by Vannitsem, S., Wilks, D. S., and Messner, J. W., pp. 91–125, Elsevier, 2018.

- 770 Schefzik, R., Thorarinsdottir, T. L., and Gneiting, T.: Uncertainty quantification in complex simulation models using ensemble copula coupling, *Statistical Science*, 28, 616–640, <https://doi.org/10.1214/13-STS443>, 2013.
- Scheuerer, M.: Probabilistic quantitative precipitation forecasting using ensemble model output statistics, *Quarterly Journal of the Royal Meteorological Society*, 140, 1086–1096, <https://doi.org/10.1002/qj.2183>, 2014.
- Scheuerer, M. and Hamill, T. M.: Variogram-Based Proper Scoring Rules for Probabilistic Forecasts of Multivariate Quantities, *Monthly Weather Review*, 143, 1321–1334, <https://doi.org/10.1175/MWR-D-14-00269.1>, 2015.
- 775 Scheuerer, M., Hamill, T. M., Whitin, B., He, M., and Henkel, A.: A method for preferential selection of dates in the Schaake shuffle approach to constructing spatio-temporal forecast fields of temperature and precipitation, *Water Resources Research*, 53, 3029–3046, <https://doi.org/10.1002/2016WR020133>, 2017.
- Schuhen, N., Thorarinsdottir, T. L., and Gneiting, T.: Ensemble model output statistics for wind vectors, *Monthly Weather Review*, 140, 3204–3219, <https://doi.org/10.1175/MWR-D-12-00028.1>, 2012.
- 780 Sklar, A.: Fonctions de répartition à  $n$  dimensions et leurs marges, *Publications de l'Institut de Statistique de l'Université de Paris*, 8, 229–231, 1959.
- Su, P.: NORTARA: Generation of Multivariate Data with Arbitrary Marginals, <https://CRAN.R-project.org/package=NORTARA>, r package version 1.0.0, 2014.
- 785 Taillardat, M., Mestre, O., Zamo, M., and Naveau, P.: Calibrated ensemble forecasts using quantile regression forests and ensemble model output statistics, *Monthly Weather Review*, 144, 2375–2393, <https://doi.org/10.1175/MWR-D-15-0260.1>, 2016.
- Thorarinsdottir, T. L. and Gneiting, T.: Probabilistic forecasts of wind speed: ensemble model output statistics by using heteroscedastic censored regression, *Journal of the Royal Statistical Society: Series A (Statistics in Society)*, 173, 371–388, <https://doi.org/10.1111/j.1467-985X.2009.00616.x>, 2010.
- 790 Thorarinsdottir, T. L., Scheuerer, M., and Heinz, C.: Assessing the calibration of high-dimensional ensemble forecasts using rank histograms, *Journal of computational and graphical statistics*, 25, 105–122, <https://doi.org/10.1080/10618600.2014.977447>, 2016.
- Van Schaeybroeck, B. and Vannitsem, S.: Ensemble post-processing using member-by-member approaches: theoretical aspects, *Quarterly Journal of the Royal Meteorological Society*, 141, 807–818, <https://doi.org/10.1002/qj.2397>, 2015.
- Vannitsem, S., Wilks, D. S., and Messner, J.: *Statistical postprocessing of ensemble forecasts*, Elsevier, 2018.
- 795 Vannitsem, S., Bremnes, J. B., Demayer, J., Evans, G. R., Flowerdew, J., Hemri, S., Lerch, S., Roberts, N., Theis, S., Atencia, A., Ben Boualègue, Z., Bhend, J., Dabernig, M., De Cruz, L., Hieta, L., Mestre, O., Moret, L., Odak Plenković, I., Schmeits, M., Taillardat, M., Van den Bergh, J., Van Schaeybroeck, B., Whan, K., and Ylhaisi, J.: *Statistical Postprocessing for Weather Forecasts – Review, Challenges and Avenues in a Big Data World*, Preprint, available at <http://arxiv.org/abs/2004.06582>, 2020.
- Wilks, D. S.: Multivariate ensemble Model Output Statistics using empirical copulas, *Quarterly Journal of the Royal Meteorological Society*, 141, 945–952, <https://doi.org/10.1002/qj.2414>, 2015.
- 800 Wilks, D. S.: On assessing calibration of multivariate ensemble forecasts, *Quarterly Journal of the Royal Meteorological Society*, 143, 164–172, <https://doi.org/10.1002/qj.2906>, 2017.
- Wilks, D. S. and Wilby, R. L.: The weather generation game: a review of stochastic weather models, *Progress in Physical Geography: Earth and Environment*, 23, 329–357, <https://doi.org/10.1177/030913339902300302>, 1999.
- 805 Williams, R. M., Ferro, C. A. T., and Kwasniok, F.: A comparison of ensemble post-processing methods for extreme events, *Quarterly Journal of the Royal Meteorological Society*, 140, 1112–1120, <https://doi.org/10.1002/qj.2198>, 2014.

Yuen, R., Baran, S., Fraley, C., Gneiting, T., Lerch, S., Scheuerer, M., and Thorarinsdottir, T.: ensembleMOS: Ensemble Model Output Statistics, <https://CRAN.R-project.org/package=ensembleMOS>, R package version 0.8.2, 2018.

Ziel, F. and Berk, K.: Multivariate Forecasting Evaluation: On Sensitive and Strictly Proper Scoring Rules, Preprint, available at <http://arxiv.org/abs/1910.07325>, 2019.

810