Consistency of the multi-model CMIP5/PMIP3-past1000 ensemble

O. Bothe^{1,2,*}, J. H. Jungclaus¹, and D. Zanchettin¹

¹Max Planck Institute for Meteorology, Bundesstr. 53, 20146 Hamburg, Germany
 ²University of Hamburg, KlimaCampus Hamburg, Hamburg, Germany
 *now at: Leibniz Institute of Atmospheric Physics at the University of Rostock, Kühlungsborn, Germany

Correspondence to: O. Bothe (ol.bothe@gmail.com)

Abstract. We present an assessment of the probabilistic and climatological consistency of the CMIP5/PMIP3 ensemble simulations for the last millennium relative to proxy-based reconstructions under the paradigm of a statistically indistinguishable ensemble. We evaluate whether simulations and reconstructions are compatible realizations of the unknown past climate evolution. A lack of consistency is diagnosed in surface air temperature data for the Pacific, European and North Atlantic regions. On the other hand, indications are found that temperature signals partially agree in the western tropical Pacific, the subtropical North Pacific and the South Atlantic. Deviations from consistency may change between sub-periods, and they may include pronounced opposite biases in different sub-periods. These distributional inconsistencies originate mainly from differences in multicentennial to millennial trends. Since the data uncertainties are only weakly constrained, the frequent over-dispersive distributional relations prevent the formal rejection of consistency of the simulation ensemble. The presented multi-model ensemble consistency assessment gives results very similar to a previously discussed single-model ensemble suggesting that structural and parametric uncertainties do not exceed forcing uncertainties.

1 Introduction

The fifth phase of the coupled model intercomparison project (CMIP5, Taylor et al., 2012) incorporates, for the first time, paleoclimate simulations in its suite of numerical experiments. The last 1000 yr of the pre-industrial period are the most recent key-period identified by the Paleoclimate Modelling Intercomparison Project Phase III (PMIP3, Braconnot et al., 2012). In contrast to the traditional time-slice simulations for specific periods of the past (e.g. Last Glacial

Maximum), the PMIP3 "past1000" experiments are transient simulations covering 850 to 1850 AD with time-varying estimates for external drivers, such as orbital, solar, volcanic and land-use climate forcings (Schmidt et al., 2011). The past1000-ensemble bridges a gap between the unperturbed control simulations and the historical simulations for the last 150 yr. It provides simulated estimates of a climate only slightly different from today. Since the ensemble allows for detailed comparisons with climate reconstructions it assists in improving our understanding of past climate forcings and naturally-forced climate variability and, in turn, in fingerprinting anthropogenic climate change (Hegerl et al., 2007; Sundberg, 2012; Schmidt et al., 2013). Assessing the quality of our simulations against paleoclimate estimates provides essential testbeds for our climate models (e.g. Schmidt et al., 2013).

Commonly, validation considers how accurately a simulated data set agrees to with the observational data in terms of matching patterns (e.g. Taylor, 2001). Comparison of simulations and Comparing simulations with reconstructions implicitly interprets both as representations of the same past. Based on this, their agreement may be taken as validation of the model and their disagreement may highlight model deficiencies. However, we have to take into account the considerable uncertainties in the reconstructions. Thus, we propose that it is appropriate in the past1000context to assess the consistency of the simulations applying methods from weather-forecast verification following, e.g. Annan and Hargreaves (2010) and Marzban et al. (2011) prior to any subjective comparison. This means that comparisons. Therefore, we have to ask whether (e.g. Hargreaves et al., 2011; Bothe et al., 2013): Do simulations and reconstructions represent compatible realizations of the unknown past climatein terms of the distribution from which its trajectory samples under the operating forcings? Answering

this question implies to establish within confidence margins whether the data samples can be assumed to stem from a common distribution whose shape is constrained by external perturbations and forcings to the climate system. Such consistency provides confidence that simulations and reconstructions indeed describe the same property and include comparable amounts of forced and internal variability. The paradigm of a statistically indistinguishable ensemble offers a theoretical basis to evaluate the consistency of simulation ensembles with reconstructions. We use this framework to assess the ensemble consistency of the past1000 multi-model-ensemble with the global temperature field reconstruction by Mann et al. (2009) and two regional areaaveraged temperature reconstructions for Central Europe (Dobrovolný et al., 2010) and Southwestern North America the North American Southwest (Wahl and Smerdon, 2012). We also discuss results for a subset of recalibrated northern hemisphere mean temperature reconstructions (Frank et al., 2010). We interpret the past1000-ensemble in terms of a probabilistic ensemble of realizations of climate variability and test whether it reliably represents the reconstructed distribution including the reconstruction uncertainty.

The current study extends our previous work (Bothe et al., 2013), which tested the consistency of the PMIP3-compliant ensemble of simulations for the last millennium performed with the COSMOS version of the MPI Earth System Model (MPI-ESM) developed at the Max Planck Institute for Meteorology (COSMOS-Mill ensemble, Jungclaus et al., 2010). The ensemble spans a number of forcing and initial conditions. We found that the COSMOS-Mill ensemble commonly lacks consistency with a set of reconstructions for Northern Hemisphere mean temperature and with the global temperature field reconstruction by Mann et al. (2009). However, its representations of Central European annual mean temperature are consistent with the reconstruction by Dobrovolný et al. (2010).

The PMIP3-past1000 multi-model-ensemble allows considering consideration of the consistency of our paleoclimate-simulations with reconstructions not only under with initial and forcing condition uncertainties (like for the COSMOS-Mill ensemble) but especially under also and especially with structural uncertainties in the models and the different parametric choices in different models and different structural uncertainties in the models (Mauritsen et al., 2012; Tebaldi and Knutti, 2007). Structural sources of uncertainties include on the most basic level different horizontal and vertical grids in all model compartments (atmosphere, ocean, land) but also the prescribed climatologies (e.g. ozone) and the formulation of, e.g. snow and ice albedo (compare, e.g. Mauritsen et al., 2012). Even if different models share certain components or portions of the same numerical code (e.g. the model genealogy of Masson and Knutti, 2011), the tuning of associated parameters likely differs. The models contributing to CMIP5 and PMIP3 (Taylor et al., 2012) generally

represent an improvement over the previous generation of models which contributed to CMIP3 (Meehl et al., 2007). In particular, some models (including MPI-ESM) provide paleo-simulations at the same resolution as the historical and the future scenario simulations. Schmidt et al. (2013) emphasise the importance of such a setup for paleo-simulations to be useful in assessing the quality of simulations of the 20th century and of future climate projections. However, in contrast to the COSMOS-Mill ensemble, none of the past1000-simulations were performed including performed include calculations of a carbon cycle. We consider the multi-model analysis to clarify the consistency of simulations under the parametric differences between the models and the common or distinct structural uncertainties of the models (e.g. Sanderson and Knutti, 2012). Therefore, we do not expect a priori increased consistency compared to our earlier results (Bothe et al., 2013).

Section 2 gives details on the methodological approach and the employed data before Sect. 3 presents we present results on the consistency of the past1000-ensemble with the reconstructions - In Sect. 4, we and identify sources for the found (lack of) consistency in Sect. 3. In Sect. 4 we discuss our results. Short concluding remarks close the manuscript.

2 Methods and data

2.1 Methods

To build confidence in a simulation ensemble we may either consider the accuracy of its members in reproducing a given (observed) target (e.g. following Taylor, 2001) or assess its statistical consistency with a target data set (see Marzban et al., 2011). The evaluation of ensemble-consistency follows the paradigm of a statistically indistinguishable ensemble (for a more detailed discussion of the methods see, e.g. Bothe et al., 2013). The underlying null-hypothesis is that the verification target and the simulations are samples from a larger common distribution and therefore exchangeable (Annan and Hargreaves, 2010; Rougier et al., 2012). In the paleoclimate context, climate reconstructions are our best estimate of an observed target.

We analyse the ensemble-consistency based on two points of view. Firstly, probabilistic consistency considers the multivariate distribution of ensemble and verification data, and, secondly, climatological consistency considers the climatological distribution of the individual simulations (e.g. Johnson and Bowler, 2009; Marzban et al., 2011; Wilks, 2011).

The probabilistic evaluation addresses how the frequencies of occurrence of the ensemble data compare to those of the verification data. It allows assessing the ensemble variance but also detecting biases. The climatological evaluation analyses the climatological variance and the biases within the individual ensemble members in relation to that of the verification data.

We can assess the probabilistic component of consistency for an ensemble by sorting and ranking the target against the ensemble data (e.g. Anderson, 1996; Jolliffe and Primo, 2008; Annan and Hargreaves, 2010; Marzban et al., 2011; Hargreaves et al., 2011). Counts of the calculated ranks are displayed as histograms following Anderson (1996). Under the null hypothesis of exchangeability (i.e. indistinguishability) of the distributions, a histogram should be flat since frequencies of observed and ensemble estimated data agree for a consistent ensemble (Murphy, 1973). The criterion of flatness does not state that the ensemble indeed is consistent (as discussed by, e.g. Hamill, 2001; Marzban et al., 2011) but it is a necessary condition for our ensemble to be a reliable representation relative to the chosen verification.

Marzban et al. (2011) emphasize the climatological consistency. They propose to evaluate it by plotting the difference between the simulated and the target quantiles against the target quantiles. For a consistent simulation, such residual quantile-quantile plots should display a flat outcome at zero. Residual quantile-quantile (r-q-q) plots ease the interpretation compared to conventional quantile-quantile plots. Marzban et al. (2011) and Bothe et al. (2013) provide more details on the advantages of the r-q-q plots.

Residual quantiles and rank counts provide an easily understandable visualization easily understandable visualizations of deviations of the ensemble relative to the verification data. In r-q-q plots, biases of the ensemble data are seen as displacements from the zero y = 0 line. A positive slope in the residual quantiles highlights an overestimation of the difference of the quantiles to the mean (i.e. the variance) compared to the target quantiles. Such a too wide data set called over dispersive, indicating an over-dispersive data set. On the other hand, a negative slope highlights an underestimation of the variance, a too narrow data set, to which we refer to as under-dispersive.

In rank histograms, dome-shapes (U-shapes) indicate too wide (too narrow) probabilistic distributions, i.e. verification data are more often close to (distant from) the mean of the distribution compared to the simulation ensemble. Positive (negative) slopes represent negative (positive) ensemblebiases, i.e. the target data over-populate high (low) ranks.

We use the χ^2 goodness-of-fit test to test for the consistency of a rank count with the uniform, i.e. flat, null hypothesis. Jolliffe and Primo (2008) provide a decomposition of the test to further consider individual deviations from the expected flat outcome. These are, among others, bias and spread deviations. Goodness-of-fit statistics are presented for these two single-deviation tests and the full test and discussed in terms of their *p* values with respect to the upper 90 % critical values (for single deviation tests it equals 2.706, for the full test with 8 degree of freedom of the χ^2 distribution it equals 13.362).

The analyses require the target to provide an accurate representation of the past climate trajectory, a condition that is hardly met in paleoclimate studies due to the associated uncertainties (e.g. Wilson et al., 2007; Bradley, 2011; Randall et al., 2007; Schmidt et al., 2011, 2013). Otherwise, we have to include an uncertainty estimate in the ensemble data (Anderson, 1996). The reconstruction targets uncertainty estimates Uncertainty estimates for the reconstruction targets are used to inflate the simulated data.

Analyses under the paradigm of statistical indistinguishability require special care if we use them in the context of paleoclimatology. Any simulated or reconstructed time series over the last 1000 yr includes components of forced and internal variability. If we assume that our estimates of past forcings are approximately correct, there should be a common forced signal in simulations and reconstructions. However, the reconstruction data uncertainties are possibly a lower bound for the disagreement between simulations and their validation target (Schmidt et al., 2013). Our analysis identifies whether the variability, forced and internaltotal variability, i.e. externally-forced and internally-generated together, originates from distributions that are similar enough to be indistinguishable. In this case, we would state that the reconstruction and the simulation ensemble are consistent. If the variability of the ensemble data deviates significantly from that of the target, our approach identifies inconsistencies. The approach can also be used to highlight in which period the long-term signals do not agree between the reconstruction and the simulation ensemble. Arising lack of consistency in terms of the distributional characteristics indicates that the simulations and the reconstruction reconstructions provide different representations of the past climate. Such deviations are informative as they suggest a need for model and reconstruction-method improvements. However, they also limit the validity of conclusions on past climate variability, the climate response to past forcings or the anthropogenic fingerprint.

The assessment of consistency reduces, in principle, the subjectivity associated to the comparison of simulations and reconstructions (compare Bothe et al., 2013). However, the large uncertainties require re-considering reconsideration of the importance of distributional deviations (compare Hargreaves et al., 2011; Annan et al., 2011). Over-dispersion does not necessarily question the overall reliability of the ensemble (see Hargreaves et al., 2011; Annan et al., 2011). On the other hand, if a simulation ensemble is found to be too narrow or biased under considering the uncertainty, further simple comparison studies between the ensemble and the reconstruction may be misleading on the considered scales. We have to consider the suggested lack of consistency in subsequent research, but we may also conclude that, under-with the present uncertainties, comparison of simulated and reconstructed estimates is not informative. Note, consistency (i.e. exchangeability) and agreement (i.e. accuracy of the temporal patterns) may differ regionally; inconsistent regions can agree in the signal and regions lacking common signals can be consistent.

2.2 Data

If we want to achieve a robust evaluation of the consistency of a simulation ensemble, we have to consider, ideally, more than one data set and more than one parameter, not least because of the prominent uncertainties in climate reconstructions. However, the global temperature field reconstruction by Mann et al. (2009) is the only data that allows for a globally coherent evaluation of a climate parameter for the last millennium. It consists of decadally smoothed data. We further employ area-averaged temperature reconstructions focused on showing two data-sets for the last 500 yr for temperature in Central Europe (Dobrovolný et al., 2010) and Southwestern North America the North American Southwest (Wahl and Smerdon, 2012). These data sets serve as examples for area-averaged reconstructions. The focus is motivated by the assumption that reconstructions and forcing data are generally more reliable on during this period (compare Bothe et al., 2013).

The past1000-simulations available from the CMIP5 database were performed with the following models (see Table 1): BCC, CCSM4 (Landrum et al., 2013), FGOALS (Zhou et al., 2011), GISS-R (two realizations, http://data.giss.nasa.gov/modelE/ar5/), IPSL-CM5A-LR (Hourdin et al., 2012) and the current generation of MPI-ESM (e.g. including an updated version of the atmospheric general circulation component ECHAM6, Giorgetta et al., 2013; Jung-claus et al., 2013). The CSIRO PMIP3-only simulation is included in our ensemble (different from Phipps et al., 2012). We do not consider the further members of the GISS-R-ensemble (R21, R22). These would over-emphasize the influence of GISS-E2-R on the ensemble measures.

We exclude the simulation with MIROC-ESM (by the Atmosphere and Ocean Research Institute at the University of Tokyo, the National Institute for Environmental Studies, and Japan Agency for Marine–Earth Science and Technology) since it shows a problematic long-term drift (A. Abe-Ouchi, personal communication, 2012). On the other hand, a simple correction is performed for the drift of the GISS-R simulations (G. A. Schmidt, personal communication, 2012, see also Schmidt et al., 2012): we subtract a lowess fit (influence of about 600 yr) to the GISS-R pre-industrial control-run (pi-Control) from the annually resolved data of interest, i.e. the grid-point data for the field evaluation and the relevant time series for the area averaged assessment.

We generally use non-overlapping decadal means for simulations and reconstructions in the commonly included period 1000–1849 CE and anomalies relative to this period. For the data from Mann et al. (2009) we choose the central date for each decade (i.e. 1844 for the 1840s) since the data are originally decadally smoothed. Results change slightly but conclusions are the same when we employ non-overlapping decadal means for this reconstruction. We also employ the three sub-periods 1000s–1270s, 1280s–1550s, 1560s–1830s to evaluate how consistency may change over time. Inclusion of climates strongly biased from our reference period (e.g. the industrial period) would complicate our assessment of consistency focused on paleoclimates. The different grids of the simulations require interpolating the data onto a common T21-grid ($\sim 5^{\circ}$).

The global field reconstruction further allows evaluating the consistency of approximations of major climate indices which are commonly interpreted to present low-dimensional descriptors of the climate system (compare, e.g. Tsonis et al., 2011; Dima and Lohmann, 2007). Our approach is, as well for these data, a step beyond the pure "by eye" approaches of reconstruction-simulation assessment. We construct two indices as field averages over the Pacific (150° E-140° W, 24-53° N) and Atlantic (74° W–0° E, 2–53° N) domains. Higher latitudes are excluded to avoid effects from sea-ice variability. Simulated indices are calculated from surface air temperatures, in contrast to the common definition via sea-surface or upper ocean temperatures. This appears justified since the reconstruction is a hybrid representation of sea-surface and near-surface air temperature (compare Brohan et al., 2006; Mann et al., 2009) with only a minority of underlying proxies being of marine origin. The indices are denoted by AMO (Atlantic Multidecadal PDO (Pacific Decadal Oscillation) and PDO (Pacific Decadal AMO (Atlantic Multidecadal Oscillation), although our definitions differ from the convention (e.g. Zanchettin et al., 2013a, and references therein). We do not preprocess the input data and do not standardize the series. The indices accumulate globally- and regionally-forced as well as potential internal signals. Our later conclusions are robust against different definitions of the regional indices.

We have to include an uncertainty estimate in our analyses by inflating the simulation ensemble (compare Anderson, 1996). These are generally randomly selected from a Gaussian distribution with zero mean. For the regional reconstructions, the standard deviations are the reported standard errors, while, for the constructed indices, we use a standard deviation of 0.2 K which approximates the estimates given for similar indices by Mann et al. (2009). For the field reconstruction, we follow Bothe et al. (2013) and take the largest standard error ($\sigma \approx 0.1729$) reported for the Northern Hemisphere mean temperature series of Mann et al. (2009) as a reasonable uncertainty estimate for the field data.

We observe that neglecting uncertainties in the reconstruction can lead to pronounced differences in our inferences about the probabilistic and climatological consistency of the ensemble. That is, the ensemble may appear under-dispersive or even consistent excluding the uncertainties, although it is found to be over-dispersive if they are considered.

Our knowledge of past forcings is rather weakly constrained as seen in comparisons of the available reconstructions for land-use, total solar irradiance and volcanic eruptions as compiled by Schmidt et al. (2011, see also discussion by Schmidt et al., 2013). For the employed simulations, differences are especially noted in the volcanic forcing, which in turn implies that they mostly influence the annual time-

scale and pentad data. We will recurreturn to this point in our discussion of origins of (lack of) consistency in Sect. 4.

Consistency in our setting depends on the reference time. Due to the way the reconstructions are produced, it is in principle advisable to <u>center centre</u> all data on the calibration period of the reconstruction (J. Smerdon, personal communication, 2012) since this time is the reference for the calculation of uncertainties. We instead <u>center centre</u> our data over the full studied period, and thereby shift the focus from the complete comparability over pre-industrial and industrial times on to the comparability of the variability over the preindustrial time only.

We stress that our results neither allow to rank the various simulators the ranking of the various simulated realisations against one another nor to decide whether individual simulations or even the ensemble mean rather than the reconstruction are more representative of past climate variability.

3 Results

3.1 Global field consistency

Figure 1 gives a first impression of the probabilistic consistency of the past1000-ensemble with the global temperature field reconstruction by Mann et al. (2009). We display the p values of tests for a uniform outcome of the rank counts at every grid-point. The goodness-of-fit test leads to the rejection of the null hypothesis of a uniform outcome at gridpoints for a p value larger than 0.9 (red in Fig. 1). Thus, the analysis shows a lack of consistency for large areas of the globe for the full period (Fig. 1a). In contrast, the European Arctic is the only spatially extended area for which rank counts deviate significantly from uniformity for an arbitrary shorter period (Fig. 1b). Possible consistency is diagnosed elsewhere. If we test for individual deviations of bias or spread, at least one of them is significant over much of the globe for both, the full and the shorter period (Fig. 1c, and d).

This general impression has to be complemented by a detailed look at individual locations to identify the character of the inconsistencies. Considering a sample of gridpoints and three different sub-periods, residual quantile distributions display various different structures when evaluating the climatological consistency of the ensemble (Fig. 2, left panels, extended sample in the Supplement). We refer to Section 2.1 for a description of how to interpret the visualisations (see also, e.g. Bothe et al., 2013). A too wide simulated distribution arises as a common feature due to very strong overestimation of cold anomalies in the simulated data and notable overestimation of positive anomalies, i.e. we see a positive slope in the residual quantiles. The most extreme over-dispersion in Fig. 2 is seen in the top left panel for 73.1E and 63.7S. Most interestingly though, the reconstructed quantile distributions frequently display consecutive shifts between the sub-periods. In turn, equivalent changes occur in the residual quantiles of the simulated data (compare, e.g. 11.2E, 2.8S). Often, but not always, the reconstructed quantile distributions shift to more negative values. This results, at some grid-points, in the following behavior of the residuals of simulated quantiles: residuals change from being negatively biased to being positively biased with an interval of nearly negligible residuals and, thus, an interval for which reconstructed and simulated quantiles appear to be consistent. That is, we see a negative offset from y = 0in the early period but a positive offset from y = 0 in the late period. Inbetween there is an interval when the residual quantiles are close to y = 0. These shifts suggest that the mean state of the reconstruction changes between the three sub-periods, and that the simulated distributions do not follow but usually feature a rather constant climatology. Patterns of residual quantiles are generally comparable between the different simulations. However, ensemble members may feature distinct residuals especially in the distributional tails. Obviously, the sample for each sub-period is small since we use non-overlapping decadal means with the full-period consisting of 85 data points.

For the same sample, rank counts confirm probabilistically the result of a generally over-dispersive ensemble by showing predominantly dome-shapes (Fig. 2, right panels), i.e. the target data occupy too often the central ranks as again best seen in the panel for 73.1E and 63.7S. Compensating discrepancies in different sub-periods can imply consistency over the full period, e.g. opposite biases cancel each other out and the rank counts are approximately uniform over the full period (see, e.g. grid-point 11E <u>3S2.8S</u>). Inconsistencies originate from different discrepancies at different locations and at the same location for different sub-periods.

Thus, according to probabilistic and climatological considerations the ensemble appears to be often over-dispersive relative to the field reconstruction. The too-dispersive ensemble-character agrees with the findings of Bothe et al. (2013) for the COSMOS-Mill-ensemble (Jungclaus et al., 2010). Furthermore, since the distributional evaluation suggests changes over time in the relation between reconstruction and the reconstruction and the simulation ensemble, we can infer that reconstruction and simulations often do not rarely represent the same climate trajectory. Neither the single-model-ensemble (COSMOS-Mill) nor the past1000 multi-model-ensemble reliably represents the climate evolution suggested by the reconstruction. However, this may be due to the uncertainties associated with the verification data.

3.2 Consistency Sources of indices disagreement for the global fields

Since the North Pacific and North Atlantic are of particular interest in assessments of low-frequency climate variability and the field evaluation indicates at best limited consistency there (compare Fig. 1), we next consider surface air temperature indices for both domains (see Sect. 2.2 for details). Accounting for uncertainty in the reconstructions, the full-period residual distributions of simulated Pacific (PDO) and Atlantic (AMO) time-series arise as to some extent over-dispersive (Fig. 5a and c) mainly due to an overestimation of the tails especially for negative anomalies. Residuals are nearly negligible for some simulations. The 90envelope for a block-bootstrap approach marginally includes the zero line of consistency for the PDO. Thus, the sampling uncertainties prevent rejecting consistency.

For both indices, the reconstructed distributions for sub-periods shift from mainly positive temperature anomalies towards negative anomalies (Fig. 5b and d). The associated residual quantiles resemble the temporal development of residual grid-point-data quantiles. A slight cold bias in the early sub-period with especially large deviations for the cold tail changes to a generally over-dispersive relation in the latest sub-periods. Distributional changes between the last two sub-periods are less prominent for PDO compared to AMO.

The full-period rank counts are significantly over-dispersive for both PDO and AMO (Fig. 5e and g) according to the goodness-of-fit test, indicating that the ensemble is not probabilistically consistent with the reconstruction for both indices. The bootstrapped intervals confirm the rejection of consistency although only marginally for the AMO.

For the sub-periods, the simulation ensembles are significantly biased for both indices in the early and significantly over-dispersive in the central sub-period (Fig. 5f and h). Bias and spread are significant for the PDO in the last sub-period, but only bias is significant for the AMO then (Fig. 5f and h). Especially prominent are the over-dispersion for the late-period PDO estimates and the bias for the late-period AMO.

Thus, the regional indices confirm the field assessment result of a simulation-ensemble that tends to be over-dispersive relative to the global field reconstruction. Again, the simulation data do not reproduce the notable ehanges in the reconstructed distributions.

3.3 Consistency of regional reconstructions

We consider additional regional area-averaged temperature reconstructions to evaluate whether the mixed result relative to the field reconstruction is representative. Already the prominent uncertainty of climate reconstructions requires such additional evaluations.

We show only results for annual central European and annual Southwestern North America temperatures starting from 1500CE. Other regional reconstructions were assessed as well but are not discussed in depth. Accounting for uncertainties in the reconstructions, residual quantile distributions indicate often full-period over-dispersion for the decadal Central European temperature data (Fig. 6a). On the other hand, the data for Southwestern North America is mainly consistent (Fig. 6b). Nevertheless deviations occur for some simulations but are not significant. These include an over- as well as an under-estimation of the cold tail and an over-estimation of the warm tail. The differences among simulations are more diverse for climatological residual quantiles relative to the Southwestern North America reconstruction compared to the results concerning the large-scale indices and the grid-point data. Climatological relations can differ remarkably for different regional reconstructions as exemplified by the Central European and Southwestern North American data.

Rank histograms (Fig. 6c and d) indicate that the ensemble is probabilistically consistent with the Southwestern North American reconstruction, but the χ^2 goodness-of-fit test leads us (only just) to reject uniformity for the European data rank counts at the considered one-sided 90level (Fig. 6c and d). Similarly, the bootstrapped intervals in Fig. 6a–d do only just result in rejecting consistency for the European data but they in principle confirm the consistency for the American data. The bootstrapped envelope also highlights the high sampling variability. We note that over-dispersive deviations are much smaller for the Central European data than for the large-scale indices or the grid-point data.

Thus, the evaluation indicates better consistency of the ensemble relative to the two semi-millennial regional annual reconstructions than for either large-scale indices or grid-point data during the full period. On the other hand, analyses on additional millennial-scale reconstructions indicate usually stronger climatologically and probabilistically over-dispersive relations with, again, notable variations in consistency over time (not shown). These regional area-average data-sets often differ more strongly in their variability from the simulation ensemble than the central European and annual Southwestern North America data.

We note that the ensemble shows negligible under-dispersion relative to the European reconstruction if we exclude the uncertainties (not shown), but it indicates slight over-dispersion under uncertainties (Fig. 6a). One could argue that, for an ideal ensemble, such rather weak opposite deviations indicate a consistent ensemble and only an over-estimation of the target uncertainties .

4 Comparison of simulations and reconstructions – sources of disagreement

Significant distributional inconsistencies possibly render moot the evaluation of Distributional inconsistencies reduce the value of assessing the agreement between simulations and reconstructions. In the best case, the associated uncertainties blur the common signals. In the worse case, the estimates do not represent the same distribution. Neverthe-

less, the mutual (dis)agreement sheds light on the shortcomings of models and reconstructions.

3.1 Global temperature fields

The lack of consistency of the past1000-ensemble is slightly less prominent compared to the COSMOS-Mill ensemble for the global temperature field for the full and the subperiods and for the full and single tests. However, Bothe et al. (2013) used interannually resolved (but decadally smoothed) data and here we use non-overlapping decadal averages. For both ensembles, deviations of the simulations from the reconstruction differ strongly between different sub-periods and even include opposite deviations. It seems that overall over-dispersion is less prominent for the decadally resolved multi-model past1000-ensemble than for the interannually resolved COSMOS-Mill ensemble. The deviations visualized We note that deviations in spread and bias are nevertheless pronounced in Fig. 2 and their similarity between individual simulations in the past1000-ensemble suggest that this is more that individual simulations show rather similar climatological residual quantiles relative to the reconstruction. This leads us to the conclusion that the less prominent over-dispersion compared to the COSMOS-Mill ensemble is not so much due to the temporal resolution than-multi-model character of the past1000 ensemble and to differences in the used forcing data sets, but that it is mostly due to the ensemble characteristics fact that the analyses of the two ensembles use different temporal resolution.

To identify sources of disagreement we first consider mapped correlation coefficients between simulations and reconstructions (Fig. 3). Again we employ non-overlapping decadal means. Each simulated and reconstructed time-series represents one realization of a climate response to the employed radiative forcing perturbations modulated by the internal variability. We also expect differences in parametrisations and methodologies to affect the outcome. We consider correlation analysis as an example for common tools a universal method in studies comparing simulations and reconstructions. Finding significant correlations between individual simulations and the reconstruction indicates that both data-sets to some extent feature a similar signal but it does not give information about the origin of the signal or whether the origin is common in both data sets.

Indeed, mapped correlation coefficients suggest various degrees of agreement between individual simulations and the reconstruction (Fig. 3). Correlations are significant nearly everywhere for the ensemble mean (two-sided 99 % level) and CCSM4 (two-sided 90 % level) but less widespread for the other simulations. Most simulations correlate significantly negative with the reconstruction at some grid-points over Antarctica. All simulations correlate significantly over the western tropical Pacific, the subtropical North Pacific and the South Atlantic. The simulations and the global reconstruction do not agree on the, possibly externally forced, phasing

of variations in Antarctica and the eastern and central tropical Pacific. We note that Mann et al. (2009) report a pronounced cold anomaly in the tropical Pacific for the Medieval Warm Period (MWP). Prominent gaps in significance are also visible for the ensemble mean and for CCSM4 over the sub-polar North Atlantic, the tropical Pacific, the Indian Ocean and central Eurasia. Similarities in the correlationpatterns may be interpreted as reflecting not only the intraensemble forcing-variability/-similarity but also the association between the models (compare Masson and Knutti, 2011).

Latitude-time plots of zonal means allow further comparison of the different data sets (Fig. 4). The reconstruction represents a near-global transition from positive anomalies in the first half (the MWP) to negative anomalies in the second half (Little Ice Age, LIA) of the considered 850 yr period (Fig. 4). The zonal means are possibly not representative in high southern latitudes due to data sparseness. The strongest warmth occurs at the beginning of the millennium. Episodic warmth interrupts the LIA during the 15th and 18th centuries and is generally confined south of 50° N.

The simulations neither capture the timing of the strongest warmth nor the near-global MWP-LIA transition. The ensemble generally displays near-stationary warm conditions. Short cold episodes related to assumed volcanic eruptions interrupt this warmth. Their timing, amplitude and spatiotemporal extent are similar in individual simulations. Weaker cold excursions reflect to some extent the variety of the employed forcings for reconstructed volcanic eruption properties (compare Schmidt et al., 2011). The ensemble mean differs most notably from the reconstruction in the lack of persistent northern hemispheric cold anomalies after about 1450 and in a stronger simulated cold signal in the 13th century. Otherwise it visually agrees well with the reconstruction.

We note that ensemble-mean correlation coefficients are often especially high (Fig. 3h) close to the proxy-locations employed by Mann et al. (2009). This implies stronger commonalities at those locations where our proxy-information about past climates are collected. That is, the similarity may allow inferring that simulations, reconstructions and underlying proxies as well as the forcing series relate to a similar underlying climate signal. Such inference is in accordance with the results of Schurer et al. (2013), Fernández-Donado et al. (2013) and Hind et al. (2012). On the other hand, the hypothesised common signal is concealed by the internal variability of the simulated climates and the additional sources of noise associated with simulations and reconstructions. We find no identifiable relation-relationship between the reconstruction and the simulations at the grid-point level.

3.1 Atlantic and Pacific Consistency of indices

Since the North Pacific and North Atlantic are of particular interest in assessments of low-frequency climate variability and the field evaluation indicates at best limited

consistency there (see Fig. 1), we next consider surface air temperature indices for both domains (see Sect. 2.2 for details). Accounting for uncertainty in the reconstructions, the full-period residual distributions of simulated Pacific (PDO) and Atlantic (AMO) time-series arise as to some extent over-dispersive (Fig. 5a and c) mainly due to an overestimation of the tails especially for negative anomalies resulting in a slight positive slope of the residuals. Residuals are nearly negligible for some simulations. The 90% envelope for a block-bootstrap approach marginally includes the zero line of consistency for the PDO. Thus, the sampling uncertainties prevent rejecting consistency.

For both indices, the reconstructed distributions for sub-periods shift from mainly positive temperature anomalies towards negative anomalies (Fig. 5b and d). The associated residual quantiles resemble the temporal development of residual grid-point-data quantiles. A slight cold bias in the early sub-period with especially large deviations for the cold tail changes to a generally over-dispersive relationship in the latest sub-periods. That is, we see a negative offset from y = 0 in the early period but a positive slope in the latest period. Distributional changes between the last two sub-periods are less prominent for PDO compared to AMO.

The full-period rank counts are significantly over-dispersive for both PDO and AMO (Fig. 5e and g) according to the goodness-of-fit test, indicating that the ensemble is not probabilistically consistent with the reconstruction for both indices because the target data over-populates the central ranks. The bootstrapped intervals confirm the rejection of consistency although only marginally for the AMO.

For the sub-periods, the simulation ensembles are significantly biased for both indices in the early and significantly over-dispersive in the central sub-period (Fig. 5f and h), i.e. we see an overpopulation of the high or low ranks in the early period but a dome-shape in the central period. Bias and spread are significant for the PDO in the last sub-period, but only bias is significant for the AMO in this period (Fig. 5f and h). Especially prominent are the over-dispersion for the late-period PDO estimates and the bias for the late-period AMO. The sub-period results are sensitive to the specific regional definitions of our indices, but the general lack of consistency is robust.

Thus, the regional indices confirm the field assessment result of a simulation-ensemble that tends to be over-dispersive relative to the global field reconstruction. Again, the simulation data does not reproduce the notable changes in the reconstructed distributions.

3.2 Sources of disagreement of Atlantic and Pacific indices

For the indices considered in the present study, the ensemble mean and the reconstruction evolve, by eye, similarly for the Atlantic and Pacific indices since about 1650 (Fig. 5i and j). Amplitudes agree less than tendencies. The most prominent example of differences in long-term trends leading to biased estimates is the different timing of medieval warmth. Note further the strong disagreement due to, on average, colder reconstructed indices from the 14th to 17th centuries for the PDO and from the 16th to 18th centuries for the AMO.

The indices display some intra-ensemble and ensemblereconstruction agreement. Again we discuss correlations as example for common practices. Ensemble-mean indices correlate at $r \approx 0.5$ with the reconstructed ones. Correlations with the reconstructed index are larger than 0.5 for the PDO in FGOALS and for the AMO in CSIRO. Correlations among simulations larger than 0.5 are only found for the AMO and most prominently for MPI-ESM and CSIRO. PDO and AMO correlate strongest in the reconstruction, CSIRO and the ensemble mean (r > 0.8). If the analysis is repeated for globally detrended data, no strong correlations are seen between the reconstructed and simulated indices.

We did not discuss regional average indices in Bothe et al. (2013), but in both regions the COSMOS-Mill simulations displayed more variability than the reconstruction and, for the North Atlantic, the ensemble-consistency changed strongly between the considered sub-periods.

3.3 Regional temperaturesConsistency of regional and hemispheric reconstructions

The assessment of the global field reconstruction raises the question whether the ensemble displays comparable lack of consistency and related sources of disagreement relative to other reconstructions. We consider additional regional area-averaged temperature reconstructions to evaluate whether the mixed result relative to the Mann et al. (2009) field reconstruction is representative. Already the prominent uncertainty of climate reconstructions requires such additional evaluations.

We show only results for annual temperatures for central Europea (Dobrovolný et al., 2010) and the North American Southwest (Wahl and Smerdon, 2012) starting from 1500 CE. Other regional reconstructions were assessed as well but are not discussed in depth. Accounting for uncertainties in the reconstructions, residual quantile distributions show often a positive slope over the full-period and therefore indicate over-dispersion for the decadal Central European temperature data (Fig. 6a). On the other hand, the data for the North America Southwest is mainly consistent (Fig. 6b) with residuals close to y = 0. Nevertheless deviations occur for some simulations but are not significant. These include an over- as well as an under-estimation of the cold tail and an over-estimation of the warm tail. The differences among simulations are more diverse for climatological residual quantiles relative to the Southwestern North America reconstruction compared to the results concerning the large-scale indices and the grid-point

data. Climatological relations can differ remarkably for different regional reconstructions as exemplified by the Central European and North American Southwest data.

Rank histograms (Fig. 6c and d) indicate that the ensemble is probabilistically consistent with the North American Southwest reconstruction, but the χ^2 goodness-of-fit test leads us to reject uniformity for the European data rank counts, since the statistics are only marginally significant at the considered one-sided 90% level (Fig. 6c and d). Similarly, the bootstrapped intervals in Fig. 6a–d do only just (i.e. marginally) result in rejecting consistency for the European data but they in principle confirm the consistency for the American data. The bootstrapped envelope also highlights the high sampling variability. We note that over-dispersive deviations are much smaller for the Central European data than for the large-scale indices or the grid-point data.

The supplement provides figures for a small selection of northern hemisphere mean reconstructions of annual-mean temperature from the recalibration ensemble by Frank et al. (2010). Frank et al. recalibrated nine northern hemispheric mean reconstructions to different periods of observational data resulting in 521 individual reconstructions.

Consistency of the simulation ensemble with the considered subset of recalibrated reconstructions is generally limited to some sub-periods, which generally differ for the individual recalibrated reconstruction series. Over-dispersion is the most common deviation but prominent biases also occur over individual sub-periods. However, we cannot overall reject consistency relative to the reconstruction by Frank et al. (2007).

The assessment of consistency of the PMIP3-past1000 ensemble relative to the members of the Frank et al. ensemble highlights an additional feature. The recalibrated series for a specific reconstruction display different variability dependent on the specific recalibration period. Therefore the simulation ensemble can be consistent with respect to a reconstruction recalibrated to a specific period while lacking consistency with the same reconstruction recalibrated to a different period. This ambiguity highlights the inherent uncertainty in our estimates for the climate of the last millennium and stresses the necessity of increasing the quality of reconstructions, of simulations and of the external-forcing estimates used for the simulations.

However, the evaluation of regional and large scale mean temperature reconstructions indicates better consistency of the ensemble relative to the two semi-millennial regional annual reconstructions than for either large-scale indices or grid-point data during the full period. On the other hand, analyses on additional regional and northern hemisphere mean millennial-scale reconstructions indicate usually stronger climatologically and probabilistically over-dispersive relations with, again, notable variations in consistency over time (not shown). These regional area-average data-sets often differ more strongly in their variability in the simulation ensemble than the central European and annual North American Southwest data.

We note that the ensemble shows negligible under-dispersion relative to the European reconstruction if we exclude the uncertainties (not shown), but it indicates slight over-dispersion under uncertainties (Fig. 6a). One could argue that, for an ideal ensemble, such rather weak opposite deviations indicate a consistent ensemble and only an over-estimation of the target uncertainties (Persson, 2011, Appendix B; see also Hargreaves et al., 2011).

3.4 Sources of regional disagreement

Figure 6 clearly displays that there is no common signal in the regional average time-series for Central Europe for individual simulations and reconstructions. This was similarly seen for the annually-resolved Central European temperature indices of the COSMOS-Mill ensemble. Obviously internal variability and methodological uncertainties dominate over the forced variability on the decadal and the inter-annual time scale for both ensembles. However, the COSMOS-Mill ensemble is consistent with the annual data of Dobrovolný et al. (2010) on the interannual time scale. Compared to the European data, the past1000-simulated Southwestern North American North American Southwest temperature series agree slightly better with the respective reconstruction for the non-overlapping decadal means. Considering the full ensemble, no common forced signal can be found. Thus we do not further comment on the accuracy of both datasets.

3.5 Further discussion

The PMIP3-past1000 ensemble and the recalibrated hemispheric mean reconstructions clearly differ in the resolved multidecadal-to-centennial and interdecadal variability. Consequently and similar to the global field data, simulations and reconstructions often differ in their long-term multicentennial trends, e.g., the passage from the MWP to the LIA and to the industrial warming period.

4 Discussions

For the field and the index data, inconsistencies over the full period are due to a generally warmer start of the millennium in the reconstruction (compare see Fig. 34). This would be mitigated for the analysis of ensemble-consistency and for the index-agreement between GISS and the reconstruction if the drift was not corrected for GISS-R (compare Schmidt et al., 2012). Decadal temperatures and their variability are more comparable in the period of the early LIA. Shifts in reconstructed quantiles towards more negative anomalies and in simulated residuals towards a more positive bias reflect the more pronounced reconstructed MWP-LIA transition and, thus, the differences in the long-term trends. Note that specific results are sensitive to the choice of the reference pe-

riod. If we align the data sets to a different common period, specific relations relationships are going to change relative to most reconstructions further highlighting the large discrepancies between the simulations and the reconstructions. The differences in estimates of cold anomaly quantiles reflect that, on the one hand, reconstructions possibly underestimate the cooling subsequent to large volcanic eruptions (Mann et al., 2012) while, on the other hand, models may be too sensitive to the subsequent radiative forcing anomaly (e.g. Anchukaitis et al., 2012).

In view of a possible impact of the choice of forcing inputs on the simulated data, one might think about partitioning the ensemble relative to the various combinations of forcings (Table 1). Only discriminating by the volcanic forcing, this results in two sub-ensembles including, respectively BCC, IPSL, CCSM4 and GISS-R25 (using the data by Gao et al., 2008) and MPI-ESM, CSIRO and GISS-R24 (using the Crowley data, see, e.g. Crowley and Unterman, 2012). Here we exclude the FGOALS data as it does not easily fit into these two categories but considers forcings as presented by Jones and Mann (2004) which are not explicitly included in the PMIP3-protocol (see Table 1 and Schmidt et al., 2011). Comparing the two ensembles we refer to the ensembles as Gao-ensemble and Crowley-ensemble.

The Crowley-ensemble generally shows a smaller response to the prescribed volcanic eruptions than the Gao-ensemble for interannual (Supp. Fig. S6) and decadal (Supp. Fig. S7) variations (compare, e.g. 1250s, we see that the effect of the different volcanic forcing data sets is comparable to the effect of different model architectures inferred from the within sub-ensemble variations (not shown). This evaluation further implicates that 1450s, and 1810s). The Gao-ensemble displays more eruptions having an influence on the northern hemisphere mean temperature than the Crowley-ensemble, which highlights the differences between the two volcanic reconstructions. The range of the ensembles, i.e. largest minus smallest temperature anomaly for each date are comparable for the majority of dates (see Supp. Fig. S8). However, the range of the Gao-ensemble is larger than for the Crowley-ensemble for a number of cases. Whereas this is mainly due to the generally much weaker response to strong volcanic eruptions in BCC, other individual simulations can differ strongly from the three other members of the Gao-ensemble at certain dates (Supp. Fig. S6-S8). Thus, the different architectures of the models can result in differences between the simulations at least as large as those induced by different volcanic forcing data sets (see Supp. Fig. 8). That is, the implementation strategy for the volcanic forcing data and the tuning of the model may influence the results as much as the choice of the forcing data (see also discussions by Fernández-Donado et al., 2013). We note that Schmidt et al. (2013) report for the GISS-R simulations and the Gao et al. (2008) data a radiative forcing twice as strong as expected from the data for the GISS-R simulations. They attribute this fact to the implementation of the volcanic forcing data in the model. We use the ensemble member GISS-R. The temperature responses to strong volcanic eruptions from the GISS-R25 simulation considered here, which employs the Gao et al. (2008) data. The simulated impacts of strong volcanic eruptions on temperature, are among the largest but not generally exceptional (not shown, compare Fig. 4in the multi-model PMIP3-past1000 ensemble but generally not exceptional (see Supp. Fig. 6).

The representation of volcanic eruptions in simulations (but also their assessment in reconstructions) is a highly controversial topic as seen in the ongoing discussions originating from Mann et al. (2012, see also Anchukaitis et al., 2012). Beyond discussions focussed on the climate of the last millennium, Driscoll et al. (2012) report an apparent lack of skill of the CMIP5 climate models in reproducing the observed response to the well constrained 20th century eruptions, possibly linked to a poor representation of relevant dynamical features. Considering the technical handling of volcanic forcing in climate simulations, Timmreck et al. (2009) and Timmreck et al. (2010) highlighted the influence of the parameterisation of aerosol size on the magnitude of the imposed radiative forcing. Zanchettin et al. (2013b) showed for a well constrained top-of-atmosphere radiative forcing, that the simulated climate response to a strong volcanic eruption can strongly vary depending on the background climate state which is defined by ongoing internal climate variability and the presence and magnitude of additional external forcings including their forcing history.

A number of possible additional and confounding factors may influence the proxies used to reconstruct the forcing data and the temperature data (e.g. precipitation, cloudiness, general circulation). Furthermore, the simulations possibly it is possible that the simulations do not fully capture the influence of , e.g. the solar forcing due to deficiencies in the representation of atmospheric chemistry and to an only partially-resolved stratosphere. These issues become especially prominent prior to 1400 (e.g. Schurer et al., 2013). Correlations may also be dominated by short-lived episodes of large forcing which are commonly featured by simulations and reconstructions (compare Schurer et al., 2013). However, Bothe et al. (2013) also showed that it is inadequate to expect a larger agreement and consistency closer to the present should not be expected. Since the zonal means suggest similarities by filtering out regional differences, one might hypothesize that the lack of common signals between reconstructions and simulations at the grid-point level is solely due to the internal and local variability masking it.

5 Summary and conclusions

The CMIP5/PMIP3-past1000-ensemble is not generally consistent with the global temperature reconstructions by Mann

et al. (2009) on a decadal time-scale. This holds for the probabilistic and the climatological assessments. Inconsistencies between reconstructions and simulations prevent reconciling both paleoclimate estimates. Our assessment of consistency over the last millennium can be biased towards being too optimistic if existing discrepancies between different multicentennial sub-periods counter-balance each other.

The simulations and the reconstruction agree least in the tropical Pacific and the sub-polar gyre region of the North Atlantic according to our evaluation, while agreement is largest in the sub-tropical Pacific and the South Atlantic. The largescale significant correlations for some individual simulations and the ensemble mean indicate that the reconstruction and the simulation ensemble possibly include a common signal.

Robust conclusions require considering more than one data set and more than one parameter due to the large uncertainties. To this regard, the ensemble is also frequently overdispersive relative to independent area-averaged regional and northern hemisphere mean temperature reconstructions. However, the ensemble is probabilistically consistent with the reconstructed annual temperatures for the Southwestern North America North America Southwest (Wahl and Smerdon, 2012).

The PMIP3-past1000 multi-model ensemble and the COSMOS-Mill single-model ensemble (Bothe et al., 2013) give very similar results with respect to their consistency, although differences exist for the diagnosed climatological deviations, which we attribute to different handling of volcanic forcing data. So, multi-model and single-model ensembles similarly lack consistency with the reconstructions. Thus, the uncertainty due to structural differences and parametrisations in the models does apparently not exceed the uncertainties associated with different forcing and initial conditions.

The PMIP3-past1000 simulation-ensemble and a selection of global and regional validation reconstruction targets are often not exchangeable climatologically and probabilistically. Therefore they should not be regarded as representing the same climate, i.e. they should not be compared under that implicit assumption.

These results imply the following:

- The ensemble may be consistent with the verification data for either the full or for sub-periods at the gridpoint level and for area-averaged data, but only few data show consistency on both time scales.
- 2. If consistency is diagnosed only for the full period, the ensemble and the reconstruction display a comparable amount of variability and a comparable climatological range over this period, but the long-term trends differ notably. We can also conclude that the variability differs between frequency bands, e.g. the reconstruction displays larger multi-centennial but smaller decadal variability and vice versa. Furthermore, analyses depending on the background climate are hampered by the lack of sub-period-consistency.

3. If, on the other hand, the data is consistent for a subperiod, analyses on the dynamics may be valid over this period, but not necessarily in other periods. Even then we have to be careful since considering a different reference period climatology may lead to different results and the background climate may influence our assessment of the dynamics.

The deviations from consistency disclose the necessity of improvements for for improvements of simulated estimates, their reconstructed forcing data and for climate-reconstructions. The large uncertainties render difficult any firm conclusions on past climate forcing and past climate variability (e.g. Hind et al., 2012). The emerging continental-scale reconstructions of the PAGES 2K consortium (e.g. Ahmed et al., 2013) offer opportunities for further application of the method to clarify the (lack of) consistency between state-of-the-art reconstructions and the current generation of simulations.

Summarizing, probabilistic and climatological evaluations indicate consistency of the PMIP3-past1000 simulationensemble with the reconstructions in some regions and over certain sub-periods, but strong biases and/or dispersiondeviations arise in other regions and periods. This dominant feature of the analysis prevents reconciling the simulated and reconstructed time-series either with respect to a common naturally forced climate signal or with respect to an estimate of internal variability.

Supplementary material related to this article is available online at: http://\@journalurl/\@pvol/\@fpage/\@pyear/\@journalnameshortlower-\@pvol-\@fpage-\@pyear-supplement.pdf.

Acknowledgements. We acknowledge comments by Eduardo Zorita and Jochem Marotzke, which lead to notable improvements on an earlier version of the manuscript as did two anonymous reviews. O. Bothe received funding by the DFG through the Cluster of Excellence CliSAP, University of Hamburg. D. Zanchettin was supported by the BMBF (research program "MiKlip", FKZ:01LP1158A). This work contributes to the MPI-M Integrated Project Millennium and CliSAP. We acknowledge the World Climate Research Programme's Working Group on Coupled Modelling, which is responsible for CMIP, and we thank the climate modeling groups (listed in Table -1 of this paper) for producing and making available their model output. For CMIP the US-U.S. Department of Energy's Program for Climate Model Diagnosis and Intercomparison provides coordinating support and led development of software infrastructure in partnership with the Global Organization for Earth System Science Portals. This work benefited from the efforts of the Paleoclimatology Program at NOAA and its archive of reconstructions of past climatic conditions and forcings.

The service charges for this open access publication have been covered by the Max Planck Society.

References

- Ahmed, M., Anchukaitis, K. J., Asrat, A., Borgaonkar, H. P., Braida, M., Buckley, B.M., Büntgen, U., Chase, B. M., Christie, D. A., Cook, E. R., Curran, M. A. J., Diaz, H. F, Esper, J., Fan, Z., Gaire, N. P., Ge, Q., Gergis, J., González-Rouco, J. F., Goosse, H., Grab, S. W., Graham, N., Graham, R., Grosjean, M., Hanhijärvi, S. T., Kaufman, D. S., Kiefer, T., Kimura, K., Korhola, A. A., Krusic, P. J., Lara, A., Lézine, A., Ljungqvist, F. C., Lorrey, A. M., Luterbacher, J., Masson-Delmotte, V., McCarroll, D., McConnell, J. R., McKay, N. P., Morales, M. S., Moy, A. D., Mulvaney, R., Mundo, I. A., Nakatsuka, T., Nash, D. J., Neukom, R., Nicholson, S. E., Oerter, H., Palmer, J. G., Phipps, S. J., Prieto, M. R., Rivera, A., Sano, M., Severi, M., Shanahan, T. M., Shao, X., Shi, F., Sigl, M., Smerdon, J. E., Solomina, O. N., Steig, E. J., Stenni, B., Thamban, M., Trouet, V., Turney, C. S., Umer, M., van Ommen, T., Verschuren, D., Viau, A. E., Villalba, R., Vinther, B. M., von Gunten, L., Wagner, S., Wahl, E. R., Wanner, H., Werner, J. P., White, J. W., Yasue, K., and Zorita, E.: Continental-scale temperature variability during the past two millennia, Nature Geoscience, 6, 339-346, doi:http://dx.doi.org/10.1038/ngeo179710.1038/ngeo1797, 2013.
- Anchukaitis, K. J., Breitenmoser, P., Briffa, K. R., Buchwal, A., Büntgen, U., Cook, E. R., D'Arrigo, R. D., Esper, J., Evans, M. N., Frank, D., Grudd, H., Gunnarson, B. E., Hughes, M. K., Kirdyanov, A. V., Korner, C., Krusic, P. J., Luckman, B., Melvin, T. M., Salzer, M. W., Shashkin, A. V., Timmreck, C., Vaganov, E. A., and Wilson, R. J. S.: Tree rings and volcanic cooling, Nat. Geosci., 5, 836–837, doi:http://dx.doi.org/10.1038/ngeo164510.1038/ngeo1645, 2012.
- Anderson, J. L.: A method for producing and evaluating probabilistic forecasts from ensemble model integrations, J. Climate, 9, 1518–1530, 1996.
- Annan, J. and Hargreaves, J.: Reliability of the CMIP3 ensemble, Geophys. Res. Lett., 37, L02703, doi:http://dx.doi.org/10.1029/2009GL04199410.1029/2009GL041994, 2010.
- Tachiiri, Annan, J., Hargreaves, J., K.: On the and observational assessment of climate model performance, Geophys. Res. Lett., 38, L24702, doi:http://dx.doi.org/10.1029/2011GL04981210.1029/2011GL049812, 2011.
- Bothe, O., Jungclaus, J. H., Zanchettin, D., and Zorita, E.: Climate of the last millennium: ensemble consistency of simulations and reconstructions, Clim. Past, 9, 1089–1110, doi:http://dx.doi.org/10.5194/cp-9-1089-201310.5194/cp-9-1089-2013, 2013.
- Braconnot, P., Harrison, S. P., Kageyama, M., Bartlein, P. J., Masson-Delmotte, V., Abe-Ouchi, A., Otto-Bliesner, B., and Zhao, Y.: Evaluation of climate models using palaeoclimatic data, Nat. Clim. Change, 2, 417–424, doi:http://dx.doi.org/10.1038/nclimate145610.1038/nclimate1456, 2012.

O. Bothe et al.: CMIP5/PMIP3-past1000 consistency

- Bradley, R. S.: High-resolution paleoclimatology, in: Dendroclimatology, Springer, 3–15, doi:http://dx.doi.org/10.1007/978-1-4020-5725-0_210.1007/978-1-4020-5725-0_2,2011.
- Brohan, P., Kennedy, J. J., Harris, I. Tett, S. F. B., and Jones, P. D.: Uncertainty estimates in regional and global observed temperature changes: a new data set from 1850, J. Geophys. Res.-Atmos., 111, D12106, doi:http://dx.doi.org/10.1029/2005JD00654810.1029/2005JD006548, 2006,
- Crowley, T. J.: Causes of climate change over the past 1000 years, Science, 289, 270–277, 2000.
- Crowley, T. J. and Unterman, M. B.: Technical details concerning development of a 1200-yr proxy index for global volcanism, Earth Syst. Sci. Data Discuss., 5, 1–28, doi:http://dx.doi.org/10.5194/essdd-5-1-201210.5194/essdd-5-1-2012, 2012.
- Crowley, T. J., Zielinski, G., Vinther, B., Udisti, R., Kreutz, K., Cole-Dai, J., and Castellano, E.: Volcanism and the little ice age, Pages News, 16, 22–23, 2008.
- Dima, M. and Lohmann, G.: A hemispheric mechanism for the Atlantic Multidecadal Oscillation, J. Climate, 20, 2706–2719, doi:http://dx.doi.org/10.1175/JCLI4174.110.1175/JCLI4174.1, 2007.
- Dobrovolný, P., Moberg, A., Brázdil, R., Pfister, C., Glaser, R., Wilson, R., Engelen, A., Limanówka, D., Kiss, A., Halíčková, M., Macková, J., Riemann, D., Luterbacher, J., and Böhm, R.: Monthly, seasonal and annual temperature reconstructions for Central Europe derived from documentary evidence and instrumental records since AD 1500, Climatic Change, 101, 69–107, doi:http://dx.doi.org/10.1007/s10584-009-9724-x10.1007/s10584-009-9724-x, 2010.
- Driscoll, S., Bozzo, A., Gray, L. J., Robock, A. and Stenchikov, G.: Coupled Model Intercomparison Project 5 (CMIP5) simulations of climate following volcanic eruptions, J. Geophys. Res., 117, D17105, doi:http://dx.doi.org/10.1029/2012JD01760710.1029/2012JD017607, 2012.
- Fernández-Donado, L., González-Rouco, J. F., Raible, C. C., Ammann, C. M., Barriopedro, D., García-Bustamante, E., Jung-claus, J. H., Lorenz, S. J., Luterbacher, J., Phipps, S. J., Servonnat, J., Swingedouw, D., Tett, S. F. B., Wagner, S., Yiou, P., and Zorita, E.: Large-scale temperature response to external forcel, ing in simulations and reconstructions of the last millennium, Clim. Past, 9, 393–421, doi:http://dx.doi.org/10.5194/cp-9-393-201310.5194/cp-9-393-2013.
- Frank, D., Esper, J., and Cook, E. R.: Adjustment for proxy number and coherence in a large-scale temperature
 reconstruction, Geophys. Res. Lett., 34, L16709, http://dx.doi.org/10.1029/2007GL030571doi:10.1029/2007GL030571, 2007.
- Frank, D. C., Esper, J., Raible, C. C., Büntgen, U., Trouet, V., Stocker, B., and Joos, F.: Ensemble reconstruction constraints on the global carbon cycle sensitivity to climate, Nature, 463, 527–530, http://dx.doi.org/10.1038/nature08769doi:10.1038/nature08769, 2010.
- Gao, C., Robock, A., and Ammann, C.: Volcanic forcing of climate over the past 1500 years: an improved ice core-based index for climate models, J. Geophys. Res.-Atmos., 113, D23111,

doi:http://dx.doi.org/10.1029/2008JD01023910.1029/2008JD010239, 2008.

- Giorgetta, M. A., Jungclaus, J., Reick, C. H., Legutke, S., Bader, J., Böttinger, M., Brovkin, V., Crueger, T., Esch, M., Fieg, K., Glushak, K., Gayler, V., Haak, H., Hollweg, H.-D., Ilyina, T., Kinne, S., Kornblueh, L., Matei, D., Mauritsen, T., Mikolajewicz, U., Mueller, W., Notz, D., Pithan, F., Raddatz, T., Rast, S., Redler, R., Roeckner, E., Schmidt, H., Schnur, R., Segschneider, J., Six, K. D., Stockhause, M., Timmreck, C., Wegner, J., Widmann, H., Wieners, K.-H., Claussen, M., Marotzke, J., and Stevens, B.: Climate and carbon cycle changes from 1850 to 2100 in MPI-ESM simulations for the coupled model intercomparison project phase 5, J. Adv. Model. Earth Syst., doi:http://dx.doi.org/10.1002/jame.2003810.1002/jame.20038, in press, 2013.
- Hamill, T. M .: Interpretation of rank histograms for verifying ensemble forecasts, Mon. Weather Rev., 129, 550-560, 2001.
- Hargreaves, J. C., Paul, A., Ohgaito, R., Abe-Ouchi, A., and Annan, J. D.: Are paleoclimate model ensembles consistent with the MARGO data synthesis?, Clim. Past, 7, 917-933, doi:http://dx.doi.org/10.5194/cp-7-917-201110.5194/cp-7-917-2011, 2011.
- Hegerl, G. C., Crowley, T. J., Allen, M., Hyde, W. T., Pollack, H. N., Smerdon, J., and Zorita, E.: Detection of human influence on a new, validated 1500-year temperature reconstruction, J. Climate, 20, 650-666, doi:http://dx.doi.org/10.1175/JCLI4011.110.1175/JCLI4011.1, 2007.
- Hind, A., Moberg, A., and Sundberg, R.: Statistical framework for evaluation of climate model simulations by use of climate proxy data from the last millennium - Part 2: A pseudo-proxy study addressing the amplitude of solar forcing, Clim. Past, 8, 1355-1365, doi:http://dx.doi.org/10.5194/cp-8-1355-201210.5194/cp-8-1355-2012, 2012.
- Hourdin, F., Foujols, M.-A., Codron, F., Guemas, V., Dufresne, J.-L., Bony, S., Denvil, S., Guez, L., Lott, F., Ghattas, J., Braconnot, P., Marti, O., Meurdesoif, Y., and Bopp, L.: Impact of the LMDZ atmospheric grid configuration on the climate and sensitivity of the IPSL-CM5A coupled model, Clim. Dynam., 40, 2167-2192, doi:http://dx.doi.org/10.1007/s00382-012-1411-310.1007/s00382-012-1411-3, 2012.
- Johnson, C. and Bowler, N.: On the reliability and calibration of ensemble forecasts, Mon. Weather Rev., 137, 1717-1720, doi:http://dx.doi.org/10.1175/2009MWR2715.110.1175/2009MWR271Medbl, G., Covey, C., Delworth, T., Latif, M., McA-2009.
- Jolliffe, I. T. Primo, C.: and Evaluating rank histograms using decompositions of the chi-square test statistic, Mon. Weather Rev., 136, 2133-2139. doi:http://dx.doi.org/10.1175/2007MWR2219.110.1175/2007MWR22191B8310.1175/BAMS-88-9-1383, 2007. 2008.
- Jones, P. Climate and Mann, M.: over past millennia, Rev. Geophys., 42, RG2002, doi:http://dx.doi.org/10.1029/2003RG00014310.1029/2003RG000143, ECMWF, Reading, UK, 2011. 2004.
- Jungclaus, J. H., Lorenz, S. J., Timmreck, C., Reick, C. H., Brovkin, V., Six, K., Segschneider, J., Giorgetta, M. A., Crowley, T. J., Pongratz, J., Krivova, N. A., Vieira, L. E., Solanki, S. K., Klocke, D., Botzet, M., Esch, M., Gayler, V., Haak, H., Raddatz, T. J., Roeckner, E., Schnur, R., Widmann, H.,

Claussen, M., Stevens, B., and Marotzke, J.: Climate and carboncycle variability over the last millennium, Clim. Past, 6, 723-737, doi:http://dx.doi.org/10.5194/cp-6-723-201010.5194/cp-6-723-2010, 2010.

- Jungclaus, J. H., Fischer, N., Haak, H., Lohmann, K., Marotzke, J., Matei, D., Mikolajewicz, U., Notz, D., and von Storch, J.: Characteristics of the ocean simulations in MPIOM, the ocean component of the MPI-Earth system model, J. Adv. Model. Earth Syst., doi:http://dx.doi.org/10.1002/jame.2002310.1002/jame.20023, in press, 2013.
- Landrum, L., Otto-Bliesner, B. L., Wahl, E. R., Conley, A., Lawrence, P. J., Rosenbloom, N., and Teng, H.: Last millennium climate and its variability in CCSM4, J. Climate, 26, 1085-1111, doi:http://dx.doi.org/10.1175/jcli-d-11-00326.110.1175/jcli-d-11-00326.1, 2013.
- Mann, M. E., Zhang, Z., Rutherford, S., Bradley, R. S., Hughes, M. K., Shindell, D., Ammann, C., Faluvegi, G., and Ni, F.: Global signatures and dynamical origins of the Little Ice Age and Medieval Climate Anomaly, Science, 326, 1256-1260, doi:http://dx.doi.org/10.1126/science.117730310.1126/science.1177303, 2009.
- Mann, M. E., Fuentes, J. D., and Rutherford, S.: Underestimation of volcanic cooling in tree-ring-based reconstructions of hemispheric temperatures, Nat. Geosci., 5, 202-205. doi:http://dx.doi.org/10.1038/ngeo139410.1038/ngeo1394. 2012.
- Marzban, C., Wang, R., Kong, F., and Leyton, S.: On the effect of correlations on rank histograms: reliability of temperature and wind speed forecasts from finescale ensemble reforecasts, Mon. Weather Rev., 139, 295-310, doi:http://dx.doi.org/10.1175/2010MWR3129.110.1175/2010MWR3129.1, 2011.
- Masson, R.: Climate D. and Knutti, model genealogy, Geophys. Res. Lett., 38, L08703, doi:http://dx.doi.org/10.1029/2011GL04686410.1029/2011GL046864, 2011.
- Mauritsen, T., Stevens, B., Roeckner, E., Crueger, T., Esch, M., Giorgetta, M., Haak, H., Jungclaus, J., Klocke, D., Matei, D., Mikoajewicz, U., Notz, D., Pincus, R., Schmidt, H., and Tomassini, L.: Tuning the climate of a global model, J. Adv. Model. Earth Syst., 4, M00A01, doi:http://dx.doi.org/10.1029/2012MS00015410.1029/2012MS000154, 2012.
- vaney, B., Mitchell, J., Stouffer, R., and Taylor, K.: The WCRP CMIP3 multi-model dataset: a new era in climate change research, B. Am. Meteorol. Soc., 88, doi:http://dx.doi.org/10.1175/BAMS-88-9-1383-1394.
 - Murphy, A. H.: A new vector partition of the probability score, J. Appl. Meteorol., 12, 595-600, 1973.
 - Persson, A.: User Guide to ECMWF forecast products, Tech. rep.,
 - Phipps, S. J., Rotstayn, L. D., Gordon, H. B., Roberts, J. L., Hirst, A. C., and Budd, W. F.: The CSIRO Mk3L climate system model version 1.0 - Part 2: Response to external forcings, Geosci. Model Dev., 5, 649-682, doi:http://dx.doi.org/10.5194/gmd-5-649-201210.5194/gmd-5-649-2012, 2012.

- Randall, D. A., Wood, R. A., Bony, S., Colman, R., Fichefet, T., Fyfe, J., Kattsov, V., Pitman, A., Shukla, J., Srinivasan, J., Stouffer, R. J., Sumi, A., and Taylor, K. E.: Climate models and their evaluation, in: Climate Change 2007: The Physical Science Basis, Contribution of Working Group I to the Fourth Assessment Report of the Intergovernmental Panel on Climate Change, edited by: Solomon, S., Qin, D., Manning, M., Chen, Z., Marquis, M., Averyt, K. B., Tignor, M., and Miller, H. L., chap. 8, Cambridge University Press, 2007.
- Rougier, J. C., Goldstein, M., and House, L.: Second-order exchangeability analysis for multi-model ensembles, submitted, 2012.
- Sanderson, Β. M. and Knutti, R.: On the interpretation of constrained climate model ensembles, Geophys. Res. Lett., 39. L16708, doi:http://dx.doi.org/10.1029/2012GL05266510.1029/2012GL052665, 2012.
- Schmidt, G. A., Jungclaus, J. H., Ammann, C. M., Bard, E., Braconnot, P., Crowley, T. J., Delaygue, G., Joos, F., Krivova, N. A., Muscheler, R., Otto-Bliesner, B. L., Pongratz, J., Shindell, D. T., Solanki, S. K., Steinhilber, F., and Vieira, L. E. A.: Climate forcing reconstructions for use in PMIP simulations of the last millennium (v1.0), Geosci. Model Dev., 4, 33-45, doi:http://dx.doi.org/10.5194/gmd-4-33-201110.5194/gmd-4-33-2011. 2011.
- Schmidt, G. A., LeGrande, A., and Healy, R.: "Interpretation of CMIP5 model ensemble", Interactive comment on "Constraining the temperature history of the past millennium using early instrumental observations" by P. Brohan et al., Clim. Past Discuss., C393-C398, 2012.
- Schmidt, G. A., Annan, J. D., Bartlein, P. J., Cook, B. I., Guilyardi, E., Hargreaves, J. C., Harrison, S. P., Kageyama, M., LeGrande, A. N., Konecky, B., Lovejoy, S., Mann, M. E., Masson-Delmotte, V., Risi, C., Thompson, D., Timmermann, A., Tremblay, L.-B., and Yiou, P.: Using paleo-climate comparisons to constrain future projections in CMIP5, Clim. Past Discuss., 9, 775-835, doi:http://dx.doi.org/10.5194/cpd-9-775-201310.5194/cpd-9-775-2013, 2013.
- Schurer, A., Hegerl, G., Mann, M. E., Tett, S. F., and Phipps, S. J.: Separating forced from chaotic climate variability over the past millennium, J. Climate, doi:http://dx.doi.org/10.1175/JCLI-D-12-00826.110.1175/JCLI-D-12-00826.1, in press, 2013.
- Steinhilber, F., Beer, J., and Fröhlich, C.: Total solar irradiance during the Holocene, Geophys. Res. Lett., 36, L19704, doi:http://dx.doi.org/10.1029/2009GL04014210.1029/2009GL040142, Dynam., 5-6, 1301-1318, doi:http://dx.doi.org/10.1007/s00382-2009.
- Sundberg, R., Moberg, A., and Hind, A.: Statistical framework for evaluation of climate model simulations by use of climate proxy data from the last millennium - Part 1: Theory, Clim. Past, 8, 1339-1353, doi:http://dx.doi.org/10.5194/cp-8-1339-201210.5194/cp-8-1339-2012, 2012.
- Taylor, K. E.: Summarizing multiple aspects of model performance in a single diagram, J. Geophys. Res.-Atmos., 106, 7183-7192, doi:http://dx.doi.org/10.1029/2000JD90071910.1029/2000JD900719, 2001.
- Taylor, K. E., Stouffer, R. J., and Meehl, G. A.: An overview of CMIP5 and the experiment design, B. Am. Meteorol. Soc., 93, 485-498, doi:http://dx.doi.org/10.1175/BAMS-D-11-00094.110.1175/BAMS-D-11-00094.1, 2012.

- Tebaldi, C. and Knutti, R.: The use of the multimodel ensemble in probabilistic climate projections, Philos. T. Roy. Soc. Α, 365. 2053-2075. doi:http://dx.doi.org/10.1098/rsta.2007.207610.1098/rsta.2007.2076, 2007.
- Timmreck, C., Lorenz, S. J., Crowley, T. J., Kinne, S., Raddatz, T. J., Thomas, M. A., and Jungclausj, J. H.: Limited temperature response to the very large AD 1258 volcanic eruption, Geophys. Res. Lett., 36, L21708, doi:http://dx.doi.org/10.1029/2009GL04008310.1029/2009GL040083, 2009.
- Timmreck, <u>C.,</u> <u>Graf</u>, <u>H.-F.</u> Lorenz, S. J., Niemeier, U., Zanchettin, D., Matei, D., Jungclaus, J. H., and Crowley, T. J.: Aerosol size confines climate response to volcanic super-eruptions, Geophys. Res. Lett., doi:http://dx.doi.org/10.1029/2010GL04546410.1029/2010GL045464, 2010.
- Tsonis, A. A., Wang, G., Swanson, K. L., Rodrigues, F. A., and da Fontura, C. L.: Community structure and dynamics in climate networks, Clim. Dynam., 37, 933-940, 2011.
- Vieira, L. E. A., Solanki, S. K., Krivova, N. A., and Usoskin, I.: Evolution of the solar irradiance during the Holocene, Astr. Astrophys., 531, A6, doi:http://dx.doi.org/10.1051/0004-6361/20101584310.1051/0004-6361/201015843, 2011.
- R. and Smerdon, J. E.: Comparative per-Wahl. E. formance of paleoclimate field index and reconstructions derived from climate proxies and noiseonly predictors, Geophys. Res. Lett., 39, L06703, doi:http://dx.doi.org/10.1029/2012GL05108610.1029/2012GL051086, 2012.
- Wilks, On reliability D.: the of the rank his-139, togram, Mon. Weather Rev., 311-316, doi:http://dx.doi.org/10.1175/2010MWR3446.110.1175/2010MWR3446.1, 2011.
- Wilson, R., D'Arrigo, R., Buckley, B., Büntgen, U., Esper, J., Frank, D., Luckman, B., Payette, S., Vose, R., and Youngblut, D.: A matter of divergence: tracking recent warming at hemispheric scales using tree ring data, J. Geophys. Res.-Atmos., 112, D17103, doi:http://dx.doi.org/10.1029/2006JD00831810.1029/2006JD008318, 2007.
- Zanchettin, D., Rubino, A., Matei, D., Bothe, O., and Jungclaus, J. H.: Multidecadal-to-centennial SST variability in the MPI-ESM simulation ensemble for the last millennium, Cilm.
- 012-1361-910.1007/s00382-012-1361-9, 2013.
- Zanchettin, D., Bothe, O., Graf, H. F., Lorenz, S. J., Luterbacher, J., Timmreck, C., and Jungclaus, J. H.: Background conditions influence the decadal climate response to strong volcanic eruptions, J. Geophys. Res. Atmos., 118, 4090-4106, doi:http://dx.doi.org/10.1002/jgrd.5022910.1002/jgrd.50229, 2013.
- Zhou, T., Li, B., Man, W., Zhang, L., and Zhang, J.: A comparison of the Medieval Warm Period, Little Ice Age and 20th century warming simulated by the FGOALS climate system model, Chinese Sci. Bull., 56, 3028doi:http://dx.doi.org/10.1007/s11434-011-4641-3041, 610.1007/s11434-011-4641-6, 2011.

Model (Acronym)	Institute	Solar	Volcanic
bcc-csm1-1 (BCC)	Beijing Climate Center, China Meteorological Administration	Vieira	Gao
CCSM4	National Center for Atmospheric Research	Vieira	Gao
CSIRO-Mk-3L-1-2 (CSIRO)	University of New South Wales	Steinhilber	Crowley (2008)
FGOALS-gl (FGOALS)	State Key Laboratory of Numerical Modeling for Atmospheric Sciences and Geophysical Fluid Dynamics, Institute of Atmospheric Physics, Chinese Academy of Sciences	Crowley (2000), Jones and Mann (2004)	Crowley (2000), Jones and Mann (2004)
GISS-E2-R (GISS-R24)	National Aeronautic and Space Administration, Goddard Institute for Space Studies	Vieira	Crowley (2008)
GISS-E2-R (GISS-R25)	National Aeronautic and Space Administration, Goddard Institute for Space Studies	Vieira	Gao
MPI-ESM-P (MPI-ESM)	Max Planck Institute for Meteorology	Vieira	Crowley (2008)
IPSL-CM5A-LR (IPSL)	Institut Pierre Simon Laplace des sciences de l'environnement	Vieira	Gao



Fig. 1. Global assessment of the goodness-of-fit test for the field data considering uncertainties in the verification target. Plotted are lower p values for tests performed at each individual grid-point. In the upper row: full χ^2 test, in the lower row: maximum of p values for single deviation tests for bias and spread. The maximum of both p values highlights grid-points where at least one test is significant. Blue smaller than 0.1, dark to light gray in steps of 0.2 within the range between 0.1 and 0.9, red larger than 0.9. Red means rejection of the uniform null hypothesis. (**a**, **c**) full period, (**b**, **d**) for the decades from the 1240s to the 1490s.



Fig. 2. Grid-point analysis of ensemble consistency for three sub-periods: 1000s-1270s, 1280s-1550s, 1560s-1830s). Left three columns: residual quantile-quantile plots for a selection of grid-points for the first (light gray), second (dark gray) and third (colored) sub-periods. Right three columns: rank histogram counts for the selection of grid-points for the three sub-periods (first to last, light to dark gray) and the full period (black, scaled to match frequencies in sub-periods). Large (small) red squares mark grid-points where spread or bias deviations are significant over the full (from left to right the first to third sub-)period. Blue squares mark deviations which are not significant. Residual quantile plots show on the x axis the quantiles of the target and on the y axis the difference between simulated and target quantiles. Consistency of the indices for the North Pacific (PDO) and North Atlantic (AMO) regions. (a d) Residual quantile plots for (a, c)

the full period and (**b**, **d**) three sub-periods (defined as for Fig. 2) of 28 records (early, light gray, middle, dark gray, late, colored). (e-h) Rank histogram counts for (e, g) the full period and (f, h) the three sub-periods (light gray to black). Numbers are the χ^2 statistics for the periods. In (f, h) numbers refer, from left to right, to the early to late sub-periods. Blue horizontal lines give the expected average count for a uniform histogram. (i, j) Time series of the indices constructed from non-overlapping decadal means. Color-code as in legend except

for shading. Shading for residual-quantiles and rank-counts (**a**, **c**, **e**, **g**) gives the 90envelope of block-bootstrapping 2000 replicates of block-length 5. Full-period residual quantile-quantile plots (left panels), rank counts (middle panels) and time series plots (right panels) for the reconstructions by (top panels) of Central European annual temperature and (bottom panels) of Southwestern North America annual temperature. For details on the representation see the caption of Fig. 5.



Fig. 3. Mapped grid-point correlation coefficients between surface air temperature series from the considered simulations and from the reconstruction. See panel titles for individual simulations. Ensemble mean in (i). Gray (black) dots mark two-sided 90 % (99 %) confidence.



Fig. 4. Time-latitude-plots of full-field zonal mean temperature anomalies with reference to the analysis period. See panel titles for individual simulations. Ensemble mean in (i) and reconstruction in (j).



Fig. 5. Consistency of the indices for the North Pacific (PDO) and North Atlantic (AMO) regions. (**a-d**) Residual quantile-quantile plots for (**a, c**) the full period and (**b, d**) three sub-periods (defined as for Fig. 2) of 28 records (early, light gray, middle, dark gray, late, colored). (**e-h**) Rank histogram counts for (**e, g**) the full period and (**f, h**) the three sub-periods (light gray to black). Numbers are the χ^2 statistics for the periods. In (**f, h**) numbers refer, from left to right, to the early to late sub-periods. Blue horizontal lines give the expected average count for a uniform histogram. (**i, j**) Time series of the indices constructed from non-overlapping decadal means. Color-code as in legend except for shading. Shading for residual-quantiles and rank-counts (**a, c, e, g**) gives the 90% envelope of block-bootstrapping 2000 replicates of block-length 5. Residual quantile plots show on the x axis the quantiles of the target and on the y axis the difference between simulated and target quantiles.



Fig. 6. Full-period residual quantile-quantile plots (left panels), rank counts (middle panels) and time series plots (right panels) for the reconstructions by (top panels) Dobrovolný et al. (2010) of Central European annual temperature and (bottom panels) Wahl and Smerdon (2012) of North American Southwest annual temperature. For details on the representation see the caption of Fig. 5.