

# MIKLIP

## A NATIONAL RESEARCH PROJECT ON DECADAL CLIMATE PREDICTION

BY JOCHEM MAROTZKE, WOLFGANG A. MÜLLER, FREJA S. E. VAMBORG, PAUL BECKER, ULRICH CUBASCH, HENDRIK FELDMANN, FRANK KASPAR, CHRISTOPH KOTTMEIER, CAMILLE MARINI, IULIA POLKOVA, KERSTIN PRÖMME, HENNING W. RUST, DETLEF STAMMER, UWE ULBRICH, CHRISTOPHER KADOW, ARMIN KÖHL, JÜRGEN KRÖGER, TIM KRUSCHKE, JOAQUIM G. PINTO, HOLGER POHLMANN, MARK REYERS, MARC SCHRÖDER, FRANK SIENZ, CLAUDIA TIMMRECK, AND MARKUS ZIESE

A German national project coordinates research on improving a global decadal climate prediction system for future operational use.

Decadal climate prediction has progressed from being an avant-garde enterprise of only a few modeling groups to the scientific mainstream within less than a decade (Smith et al. 2007; Keenlyside et al. 2008; Pohlmann et al. 2009; Mochizuki et al. 2010; Kirtman et al. 2013; Meehl et al. 2014). Responding to both the new research opportunities and the enhanced societal requirements for information about near-term future climate change (e.g., WMO 2011; Kirtman et al. 2013), the German Federal Ministry for Education and Research has for the period 2011–19 funded a comprehensive national project on decadal climate prediction, *Mittelfristige Klimaprognose (MiKlip; midterm climate forecast)*. This paper summarizes the scientific, strategic, and structural lessons learned from MiKlip so far.

A decadal prediction system simulates not only the climate response to future natural and anthropogenic forcing but also the future evolution of internal climate variability, caused by chaotic processes. Because chaos fundamentally limits climate predictability, a decadal prediction must be initialized from the observed state of those components of the climate system that provide a multiyear “memory,” usually but not exclusively the ocean (e.g., Bellucci et al. 2015a). Relevant ocean memory arises from

the persistence of ocean heat content anomalies, especially where the atmosphere interacts with deep oceanic mixed layers, such as in the North Atlantic and North Pacific Subpolar Gyres (e.g., Mochizuki et al. 2010; Guemas et al. 2012; Matei et al. 2012b). Ocean memory possibly also arises from properly initialized ocean circulation and hence “slow” ocean dynamics [e.g., Matei et al. (2012b); a comprehensive review of the principles behind decadal prediction was recently provided by Kirtman et al. (2013)].

The quality of a decadal prediction system is assessed—in analogy to a seasonal prediction system—by performing a set of hindcasts (retrospective predictions) and by evaluating these hindcasts against the observed climate evolution. This evaluation step requires a sufficiently powerful observing system and is therefore usually limited to the period since around 1960. Assessing the gain in prediction skill that is obtained through the initialization is a core element of decadal prediction research, although for the users of such a prediction it matters little whether skill arises from the expected change in forcing or from the initialized internal variability.

The MiKlip project aims to establish and improve a decadal climate prediction system that by the end of the project can be transferred to the German

meteorological service (Deutscher Wetterdienst; DWD) for operational use. To serve this dual purpose—preoperational predictions combined with research progress—MiKlip is organized around a hub consisting of a global climate prediction system, in turn comprising the Max Planck Institute Earth System Model (MPI-ESM; Giorgetta et al. 2013) together with an initialization procedure. Around this hub, the research is organized in four modules focusing on initialization, evaluation, processes and modeling, and regionalization.

The MiKlip hub furthermore provides a central evaluation system. The evaluation system, the necessary observational data, and the entire set of MiKlip prediction results conform to the CMIP5 data standards (Taylor et al. 2012) and reside on a dedicated data server. The MiKlip server makes the prediction results and evaluation system immediately accessible to the entire MiKlip community, thereby providing a crucial interface between production on the one hand and research and evaluation on the other hand.

The structure of MiKlip differs notably from other community efforts in decadal climate prediction, especially the decadal prediction portion of phase 5 of the Coupled Model Intercomparison Project (CMIP5; see Kirtman et al. 2013; Meehl et al. 2014). CMIP5 comprises 16 different decadal prediction systems

and thus offers a much richer spectrum of modeling approaches than does MiKlip, which focuses on a single global prediction system. On the other hand, MiKlip can produce quick and tailored research responses that help modify its prediction system. MiKlip could hence cycle through a greater number of generations of its prediction system, compared to the cycle defined by the different phases of CMIP; this faster cycle enables faster learning from successive generations (see “Three generations of the global prediction system” section).

A project that conceptually rests in between MiKlip and CMIP is Seasonal-to-Decadal Climate Prediction for the Improvement of European Climate Services (SPECS; [www.specs-fp7.eu/](http://www.specs-fp7.eu/)), funded by the European Union Framework Program 7. SPECS comprises six European climate prediction systems and thus shares with CMIP the multimodel approach. SPECS shares with MiKlip the strategy to coordinate research within the project and to coordinate improvements of the prediction systems; however, SPECS is not designed to provide the same interactive cycle of prediction system improvements as MiKlip does. Overall, the approaches by MiKlip, SPECS, and CMIP complement each other.

The remainder of this paper is dedicated to the following scientific and strategic topics. The “Three generations of the global prediction system” section documents how we explored a variety of initialization methods and developed a strategy for deciding among them. These decisions have resulted in the succession of three generations of the MiKlip global decadal prediction system. The “Evaluation of prediction system generations” section demonstrates that the systematic effort in prediction evaluation and verification has led to identification of prediction skill in many new quantities, such as multiyear-mean seasonal surface temperature over Europe, Northern Hemisphere midlatitude storm tracks, the quasi-biennial oscillation (QBO), and carbon uptake by the North Atlantic. The “Processes and model development” section presents aspects of enhanced process understanding and, in particular, how the development of a volcano code package enables us to include in future predictions the occurrence of a major volcanic eruption. The “Downscaling the decadal prediction” section discusses how the regionalization of the predictions has made possible the identification of regional forecast skill. The “Discussion and conclusions” section provides a synthesis of the lessons learned from MiKlip so far.

### THREE GENERATIONS OF THE GLOBAL PREDICTION SYSTEM.

The MiKlip funding

**AFFILIATIONS:** MAROTZKE, MÜLLER, VAMBORG, KRÖGER, POHLMANN, SIENZ, AND TIMMRECK—Max Planck Institute for Meteorology, Hamburg, Germany; BECKER, KASPAR, SCHRÖDER, AND ZIESE—Deutscher Wetterdienst, Offenbach, Germany; CUBASCH, PRÖMMEL, RUST, ULBRICH, AND KADOW—Institute of Meteorology, Freie Universität Berlin, Berlin, Germany; FELDMANN AND KOTTMEIER—Institute for Meteorology and Climate Research (IMK-TRO), Karlsruhe Institute of Technology, Karlsruhe, Germany; MARINI, POLKOVA, STAMMER, AND KÖHL—Institute of Oceanography, Center for Earth System Research and Sustainability (CEN), University of Hamburg, Hamburg, Germany; KRUSCHKE—Institute of Meteorology, Freie Universität Berlin, Berlin, and GEOMAR Helmholtz Centre for Ocean Research Kiel, Kiel, Germany; PINTO—Institute for Geophysics and Meteorology, University of Cologne, Cologne, Germany, and Department of Meteorology, University of Reading, Reading, United Kingdom; REYERS—Institute for Geophysics and Meteorology, University of Cologne, Cologne, Germany

**CORRESPONDING AUTHOR:** Jochem Marotzke, Max Planck Institute for Meteorology, Bundesstrasse 53, 20146 Hamburg, Germany  
E-mail: jochem.marotzke@mpimet.mpg.de

*The abstract for this article can be found in this issue, following the table of contents.*

DOI:10.1175/BAMS-D-15-00184.1

In final form 10 June 2016  
©2016 American Meteorological Society

period is subdivided into five development stages of usually 18 months in length. Each transition from one development stage to the next marks a well-defined and easy-to-communicate point in time for collecting, synthesizing, and implementing recommendations for changes in the global prediction system. Three generations of the prediction system are now available, termed baseline0, baseline1, and prototype (Table 1). Because of the relative timing of CMIP5 and the MiKlip start, we could use the CMIP5 initialized simulations (hindcasts) as our starting point, a set that we redubbed for MiKlip use as baseline0. Already during development stage 1, we defined and performed the next set of hindcasts (baseline1), using an initialization procedure and initialization data different from baseline0. Based on the research during development stage 1, we have defined and executed during development stage 2 the experiments with the prototype system. We have not defined a prediction generation for development stage 3 (see “Discussion and conclusions” section); at this writing, we are at the beginning of development stage 4.

*From baseline0 to baseline1.* Our design of baseline1 started from the recognition that baseline0 performed poorly in the tropics. Following Matei et al. (2012b), the initial conditions in baseline0 were constructed from a simulation with the Max Planck Institute Ocean Model (MPIOM; Jungclaus et al. 2013) forced

by the National Centers for Environmental Prediction (NCEP)–National Center for Atmospheric Research (NCAR) reanalysis (Kalnay et al. 1996). The three-dimensional ocean temperature and salinity anomalies of the forced ocean run were added to the coupled model climatology; in a step with the coupled model called the assimilation run, the ocean hydrography was nudged to this sum of fields. The coupled model state resulting from the assimilation run was used as the initial condition for the 10-yr-long hindcast simulations. While this simple initialization gave excellent hindcast skill for North Atlantic sea surface temperature (SST) and even some skill in central European summer surface air temperature (Müller et al. 2012), the initialization led to degraded performance for SST in the tropics, compared to the uninitialized (historical) CMIP5 simulations (Figs. 1a,d; Müller et al. 2012; Bellucci et al. 2015b). This poor performance in the tropics may have arisen from the very simple initialization procedure, leading to a lack of balance between zonal wind stress and ocean surface pressure gradient in the coupled model (Thoma et al. 2015) or from the observations used in the procedure (e.g., McGregor et al. 2012; Lee et al. 2013; Pohlmann et al. 2016, manuscript submitted to *Geophys. Res. Lett.*).

A test suite of three-member hindcast ensembles with yearly start dates from 1961 onward explored various alternative initialization procedures. For each initialization, hindcast skill was evaluated for some

**TABLE 1. Experiments performed in MiKlip. In MPI-ESM-LR, LR stands for low resolution, T63 horizontally with 47 levels in the atmosphere and nominally 1.5° horizontal resolution and 40 levels in the ocean. In MPI-ESM-MR, MR stands for mixed resolution, T63 with 95 levels in the atmosphere and 0.4° horizontal resolution with 40 levels in the ocean.**

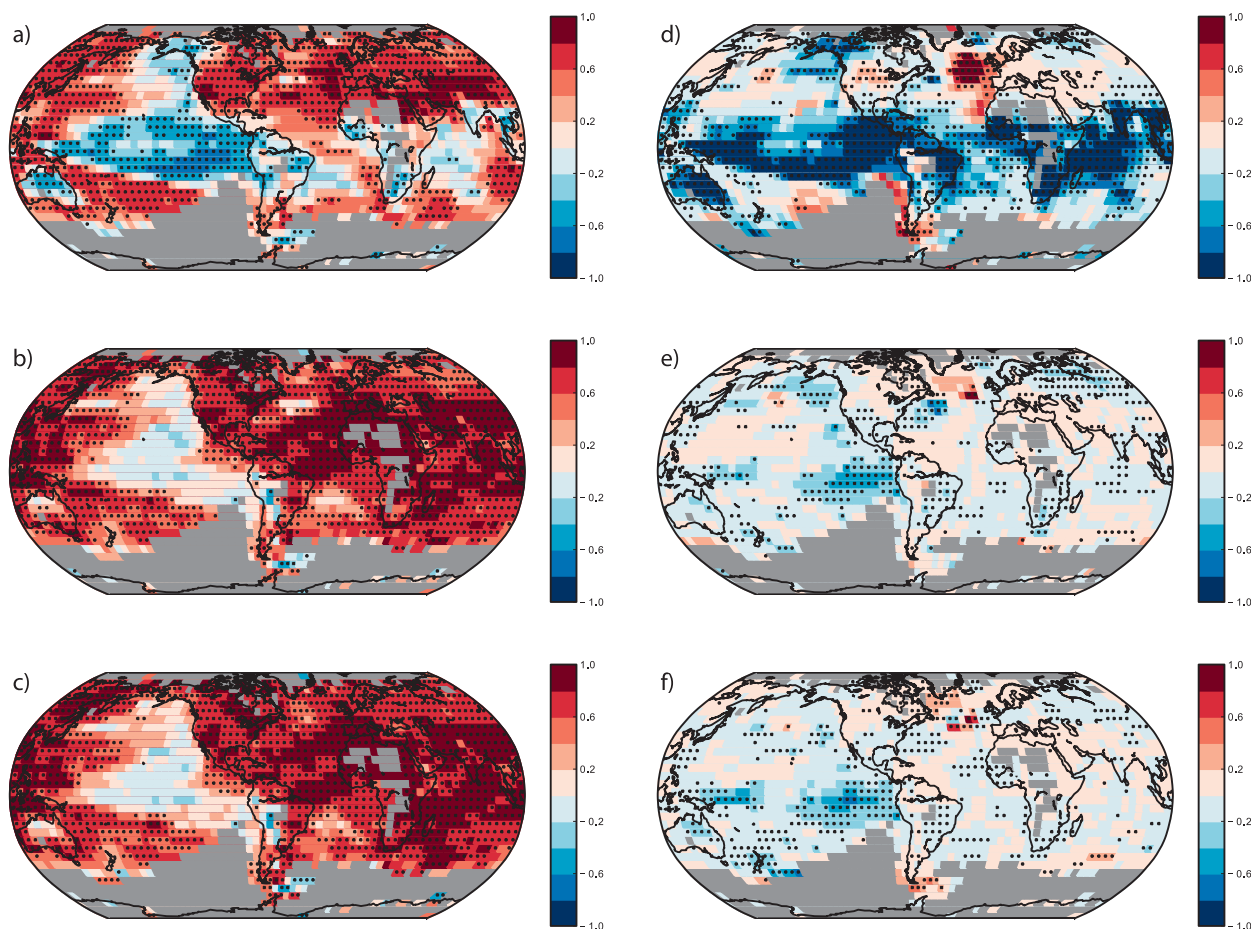
	<b>Baseline0</b>	<b>Baseline1</b>	<b>Prototype</b>
<b>Models</b>	MPI-ESM-LR MPI-ESM-MR	MPI-ESM-LR MPI-ESM-MR	MPI-ESM-LR
<b>Initialization ocean</b>	3D temperature ( <i>T</i> )–salinity ( <i>S</i> ) anomalies from MPIOM forced with NCEP–NCAR reanalysis	3D <i>T</i> – <i>S</i> anomalies from ORAS4	3D <i>T</i> – <i>S</i> (full field) from GECCO2 and from ORAS4
<b>Initialization atmosphere</b>	Assimilation run	ERA-40 and ERA-Interim; vorticity, divergence, log ( <i>p</i> ), <i>T</i> ; full field	ERA-40 and ERA-Interim; vorticity, divergence, log ( <i>p</i> ), <i>T</i> ; full field
<b>Ensemble size</b>	LR: 3 (10) MR: 3	LR: 10 MR: 5	30 (15 each with initialization from GECCO2 and ORAS4)
<b>Start years</b>	LR: 1961–2013; yearly for 3 realizations 1961–2000: five yearly for 10 realizations 2001–13: yearly for 10 realizations MR: 1961–2000: five yearly 2001–12: yearly	LR: 1961–2014: yearly MR: 1961–2013: yearly	1961–2014: yearly

predefined measures such as global-mean surface temperature, North Atlantic SST index, and, for years 2004–10, the Atlantic meridional overturning circulation (AMOC) at 26.5°N. These evaluations suggested initializing the ocean with temperature and salinity anomalies from the Ocean Reanalysis System 4 (ORAS4; Balmeda et al. 2013) reanalysis and the atmosphere from the 40-yr European Centre for Medium-Range Weather Forecasts (ECMWF) Re-Analysis (ERA-40; Uppala et al. 2005) and ECMWF interim reanalysis (ERA-Interim; Dee et al. 2011; Table 1).

Baseline1 shows much improved correlation skill for tropical surface temperature, compared to baseline0, while maintaining positive skill in North Atlantic surface temperature (Fig. 1; see also Pohlmann et al.

2013). Almost all regions with negative correlation in baseline0 show positive correlation in baseline1 (tropical Atlantic, Africa, Indian Ocean, and western Pacific). Only the eastern Pacific continues to show negative skill, although less pronounced than in baseline0, in a pattern resembling the Pacific decadal oscillation (see also Mochizuki et al. 2010; Guemas et al. 2012). The improvement in tropical SST hindcast skill in baseline1 has led to a substantial improvement also in hindcast skill for global-mean surface temperature (Pohlmann et al. 2013).

Compared against the uninitialized (historical) simulations, initialization continues to provide additional skill primarily in the North Atlantic, owing to the deep mixed layers and associated long-lived heat



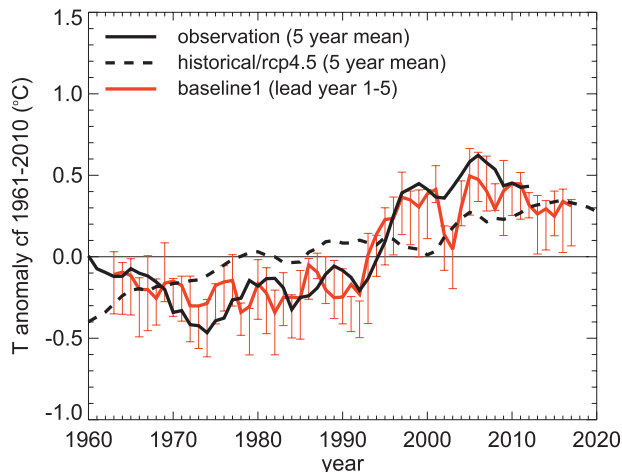
**FIG. 1.** Evolution during the MiKlip project of the ensemble-mean hindcast skill [(left) anomaly correlation; (right) with anomaly correlation of historical simulations subtracted] of surface air temperature averaged over the lead years 2–5 in the low-resolution model version MPI-ESM-LR. Observations are from Hadley Centre/Climatic Research Unit, version 4 (HadCRUT4; Jones et al. 2012). The period is 1961–2012; (a),(d) baseline0, (b),(e) baseline1, and (c),(f) prototype. Hindcast ensemble size is 3 for baseline0, 10 for baseline1, and 30 for prototype; historical ensemble size is 3 for baseline0, 10 for baseline1, and 15 for prototype. Crosses denote skill different from zero exceeding the 95% confidence level; significant negative skill indicates where the initialization causes skill degradation. The figure was created with the evaluation tool provided by the central evaluation system (Illing et al. 2014).

content anomalies there (Fig. 1e). Because the skill enhancement in the North Atlantic is supported by robust physical understanding (e.g., Matei et al. 2012b), we have confidence in this result, although the region covers only a small portion of the globe. Notice that northeastern North Atlantic SST skill relative to the historical simulations in baseline0 is inflated because of one particularly improbable historical realization within the small ensemble of three; the larger ensemble size in baseline1, both in initialized and historical simulations, means that skill assessment is more robust (see “Evaluation of prediction system generations” sections). The baseline1 hindcasts track the observed time series of North Atlantic Subpolar Gyre SST quite well and much better than the historical simulations, with the exception of a large and unexplained drop centered around year 2002 (Fig. 2). In particular, the hindcasts also show the downward trend beginning in 2005 [as was found earlier by Hermanson et al. (2014) with the Met Office decadal prediction system], and our predictions suggest that this downward trend is not reversed until the end of the current decade.

*From baseline1 to prototype.* The design of the prototype system was based on a far more comprehensive assessment compared to the design of baseline1. Suggestions for modifications were collected from each MiKlip subproject; a number of suggestions for modified initialization could readily be implemented and tested.

The first suggestion is based on the recognition that the German contribution to Estimating the Circulation and Climate of the Ocean 2 (GECCO2) ocean reanalysis (Köhl 2015) provides an improved initial state compared to its predecessor GECCO [which was used earlier in Pohlmann et al. (2009), Matei et al. (2012b), and Kröger et al. (2012)]. The model comprises higher horizontal and vertical resolution, the domain is now fully global including the Arctic, and the simulation has been extended into the most recent years. Benefits of the new assimilation can be seen in several GECCO2 solution properties crucial for decadal prediction, such as ocean heat content, which, compared to the reference simulation (without assimilation), shows reduced and more realistic interdecadal variability. The AMOC at 26.5°N agrees excellently between the reanalysis and the observations (Fig. 3; Köhl 2015).

The workflow for producing initial conditions from GECCO2 has been modified so that the data needed for the initialization are available for quasi-operational use. Such availability, ideally with no more than a 1-month delay, cannot currently be obtained through the full-blown and computationally



**FIG. 2.** SST index for the North Atlantic Subpolar Gyre (40°–60°N, 0°–60°W) from 1960 to 2020. Shown is the 5-yr running mean for the observations [Hadley Centre Sea Ice and Sea Surface Temperature dataset 1.1 (HadISST1.1); Rayner et al. 2003; solid black] and the ensemble mean over 10 realizations of MPI-ESM-LR historical simulations extended with the representative concentration pathway (RCP4.5) scenario (dashed black). Further shown is the time mean over the first five prediction years for each start year, for the ensemble mean of baseline1 MPI-ESM-LR hindcasts and predictions (solid red; last start year 2015). The whiskers show the range (minimum to maximum) of the ensemble.

intensive four-dimensional variational data assimilation (4D-Var) method on which GECCO2 is based. This drawback is overcome here by performing shorter independent optimization runs toward the end of the assimilation window and further by appending a brief unconstrained run with unadjusted forcing for the final period. This modification in the workflow might make 4D-Var more broadly applicable not only for reanalyses but also for predictions.

The second suggestion for modified initialization concerns the use of full-field rather than anomaly initialization in the ocean, reflecting a more general tendency in the decadal prediction field (Smith et al. 2013a; Meehl et al. 2014; Polkova et al. 2014). A simulation closer to the observed mean state, instead of the coupled model’s, offers conceptual advantages because some important climate processes such as sea ice formation and melt and atmospheric tropical stability are sensitive to the background state. Moreover, full-field initialization obviates the need to compute anomalies separately.

A suite of three-member test hindcast ensembles, using each of ORAS4 and GECCO2 in both anomaly and full-field ocean initialization, suggested that all three initialization alternatives to the baseline1

initialization (cf. Figs. 1b,e) led to improvements in the eastern tropical Pacific, the Indian Ocean, and the region in the northwestern North Atlantic where the three-member subensemble of baseline1 showed a relative minimum in skill (not shown). Although the skill was not improved everywhere, we concluded from the results of the initialization module (Polkova et al. 2014) and our additional test ensemble that the prototype system should use full-field initialization. The differences between ORAS4 and GECCO2 were only slight (not shown), so we used both initialization fields side by side.

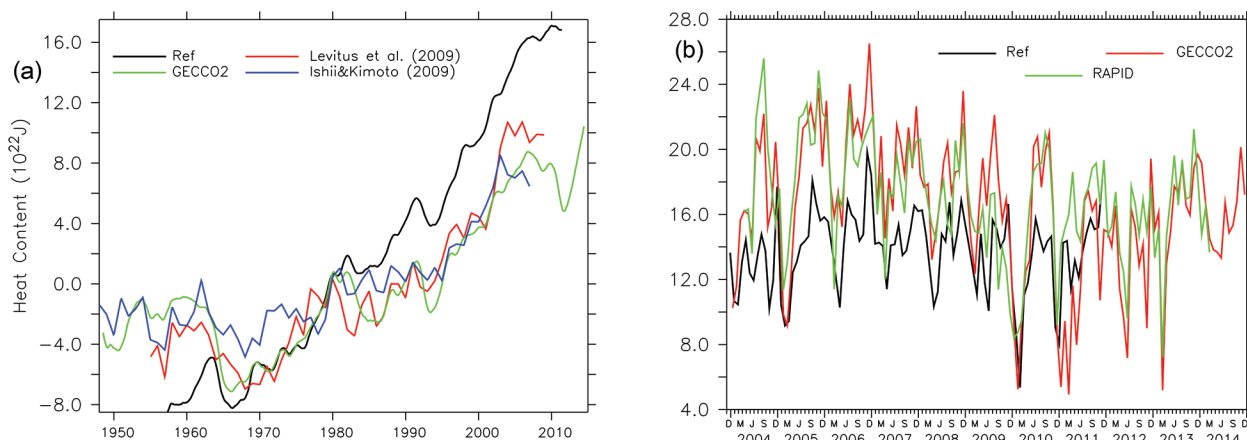
Most baseline0 and baseline1 hindcasts were performed with the Max Planck Institute Earth System Model, low resolution (MPI-ESM-LR; T63 with 47 levels in the atmosphere and nominally 1.5° horizontal resolution and 40 levels in the ocean). The Max Planck Institute Earth System Model, mixed resolution (MPI-ESM-MR; T63 with 95 levels in the atmosphere; 0.4° horizontal resolution with 40 levels in the ocean), has yielded only modest benefit in the hindcasts (Pohlmann et al. 2013), just as in the CMIP5 historical simulations (Jungclaus et al. 2013). Clear exceptions exist where use of the higher vertical resolution is essential, such as for the QBO (Pohlmann et al. 2013; see “Evaluation of prediction system generations” section). But given the computational constraints, we decided against the use of MPI-ESM-MR in the prototype system.

Instead, the prototype system employs a much larger ensemble than before. With increasing ensemble size, the ensemble-mean correlation with observations is expected to increase, while

the uncertainty of the skill estimate and the risk of finding spurious skill are expected to decrease (Murphy 1990; Kumar et al. 2001; Scaife et al. 2014a). These expectations are confirmed in baseline1 for the North Atlantic SST index and central European summer surface temperature (Fig. 4; Sienz et al. 2016). The prototype system thus comprises 30 ensemble members instead of 10, with 15 members each based on ORAS4 and GECCO2 (Table 1).

Hindcast ensembles are generated in baseline0 and baseline1 through lagged initialization, meaning that the model initial state at the nominal start day (1 January of any given start year) is taken from the state a few days earlier or later. The chaotic nature of the atmospheric model solution implies that the realizations soon drift away from each other and develop their own weather histories. But this procedure does not explore the possible ocean initial conditions that within uncertainty bounds are consistent with the available observations. Therefore, MiKlip aims at the development of alternative ensemble-generation procedures that explore the possible initial states more fully (see also Du et al. 2012).

Four procedures have been tested: empirical oceanic singular vectors (Molteni et al. 1996; Marini et al. 2016), the anomaly transform (Wei et al. 2006; Romanova and Hense 2015), a multiassimilation run approach in which the assimilation is based on several realizations of a historical run (Keenlyside et al. 2008), and the singular evolutive interpolated Kalman (SEIK) filter (Pham et al. 1998; Brune et al. 2015). Unfortunately, no robust improvement compared to the lagged initialization has been found; if there



**FIG. 3.** (a) Anomalies of the global-mean upper-ocean (0–700 m) heat content from the reference run (Ref, no assimilation) and GECCO2 in comparison to the estimates from Levitus et al. (2009) and Ishii and Kimoto (2009). (b) Comparison of the monthly mean meridional overturning at 26.5°N in Sv ( $1 \text{ Sv} = 10^6 \text{ m}^3 \text{ s}^{-1}$ ) from GECCO2 (red) and the reference run (black) with observations from Rapid Climate Change (RAPID; green; e.g., Cunningham et al. 2007; McCarthy et al. 2012). Updated from Köhl (2015), including a correction to the curve from the reference run in (a).

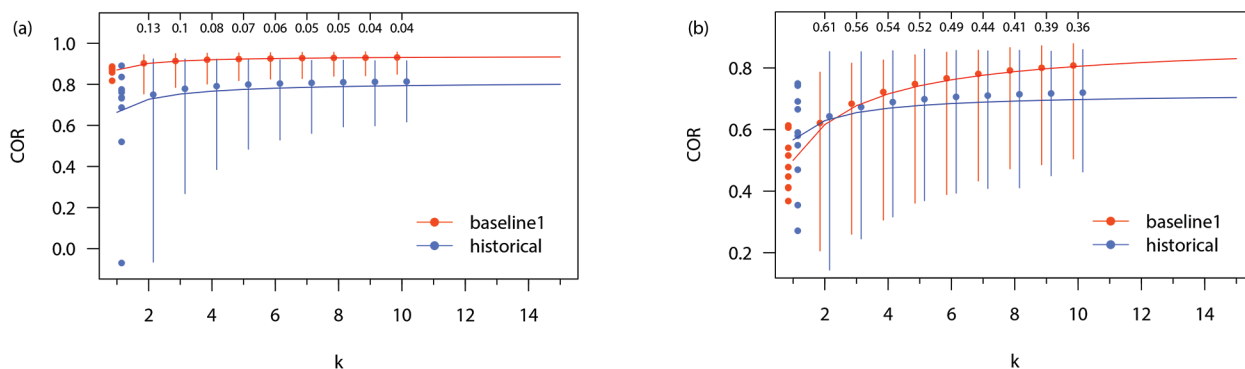
is improvement, this is compensated by additional problems such as an overestimation of the internal variability by the ensemble spread in some, though not all, variables (Marini et al. 2016). A speculative interpretation of this result suggests that on the time scales relevant here, variability even in the ocean interior might be dominated by the forcing from atmospheric internal variability. Because the more sophisticated ensemble-generation methods do not yet provide a clear path forward, we use the same lagged initialization procedure in the prototype system as in baseline0 and baseline1.

Given the large effort that went into designing and executing the prototype system, the comparison against baseline1 for surface temperature averaged over lead years 2–5 is a little sobering. We see incremental improvement in the correlation with observations, such as in the eastern tropical Pacific and the central North Atlantic (Figs. 1b,c), but the skill improvement by initialization has not increased against baseline1, except around Drake Passage and the Indian Ocean portion of the Southern Ocean (Figs. 1e,f). The anticipated improvements from the combination of enhanced ensemble size and full-field initialization have thus not materialized for all quantities.

**EVALUATION OF PREDICTION SYSTEM GENERATIONS.** The evaluation module pursues two related but distinct objectives; first, data-oriented evaluation of the prediction system and, second, process-oriented evaluation beyond the estimation

of forecast skill for standard model output. Much of the data-oriented work stems from the recognition that observational datasets often provide insufficient spatiotemporal coverage or quality to enable a comprehensive evaluation of the prediction system. Therefore, considerable work is required on these observational datasets themselves. For example, global precipitation data over both land and ocean have been reprocessed for the period 1988–2008 to deliver daily maps with a grid resolution of  $1^\circ \times 1^\circ$  and  $2.5^\circ \times 2.5^\circ$ , with a traceable estimate of the uncertainty (Schamm et al. 2014; Andersson et al. 2016a,b). As another example, variations in terrestrial water storage since 2002 have been inferred from GRACE satellite gravity measurements and used for the evaluation of the MiKlip hindcasts (Zhang et al. 2015).

The work on verification and process-oriented evaluation takes as its starting point the recommendations by Goddard et al. (2013). These include bias adjustment, typical spatial and temporal scales of aggregation, and verification of the hindcast ensemble proceeding along two lines. The first line of verification focuses on the mean square error skill score (MSESS), which tests whether the ensemble mean of a prediction outperforms a reference prediction, measured against a verification dataset. In the simple case of climatology as reference forecast, the MSESS combines the correlation between anomalies, the conditional bias (the prediction system systematically overestimates or underestimates the magnitude of anomalies), and the unconditional bias (difference



**FIG. 4.** Correlation with observations for (a) annual-mean North Atlantic ( $20^\circ$ – $60^\circ$ N,  $10^\circ$ – $80^\circ$ W) SST and (b) central European ( $40^\circ$ – $45.5^\circ$ N,  $10^\circ$ – $30^\circ$ E) summer (Jun–Aug) surface temperature for baseline1 hindcasts averaged over lead years 2–5 (red) and historical runs (blue). Shown is the dependence of the correlation on the ensemble size  $k$ ; the vertical lines are 95% confidence intervals. The dots at  $k = 1$  give the correlations for the single members. Numbers at the top are  $p$  values for the correlation skill scores of baseline1, with historical simulations as the reference prediction. The approximate theoretical correlation–ensemble size relation is given by the red (baseline1) and blue (historical) solid lines, based on Murphy (1990). The observations used are HadISST1.1 (Rayner et al. 2003) for SST and Climate Research Unit Time Series version 3.10 (CRU TS3.10) (Harris et al. 2014) for surface temperature over land. From Sienz et al. (2016), reproduced with permission ([www.schweizerbart.de](http://www.schweizerbart.de)).

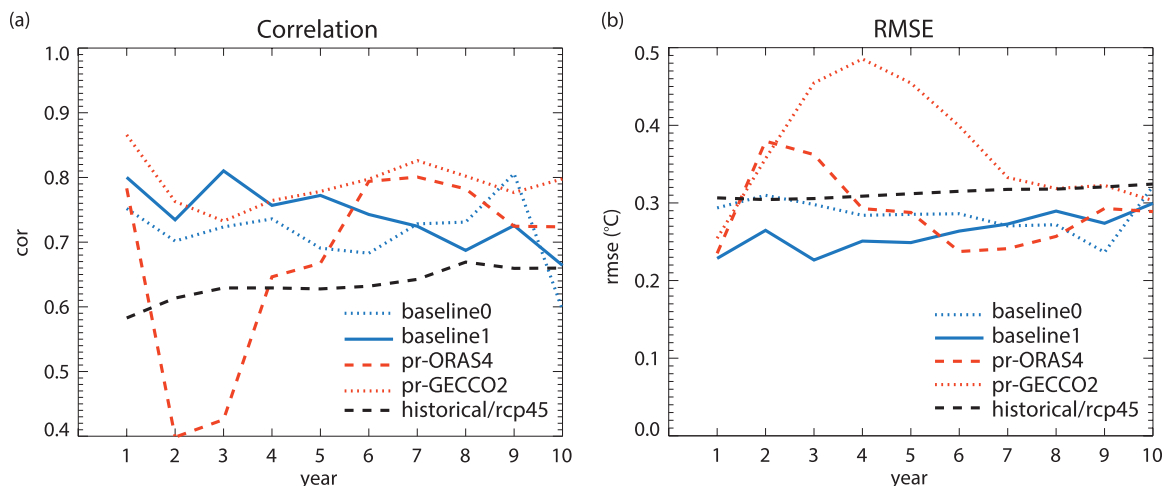
between time averages; Murphy 1988). In some results shown here, the anomaly correlation is used because the conditional bias is assumed small and the unconditional bias has been subtracted. The second line of verification focuses on the full probabilistic hindcast derived from the ensemble. We use a variant of the rank probability skill score (RPSS), which assesses whether the ensemble spread of predictions accurately represents the forecast uncertainty (e.g., Kadow et al. 2015).

The central evaluation system is constantly expanded with contributions from the MiKlip evaluation module and, together with its reference data pool for verification, resides on the same data server as the entire MiKlip prediction output. The analyses are collected into a database ensuring reproducibility and transparency. Providing the central evaluation system to the entire MiKlip project is also an effective training tool, especially for those researchers who have only recently joined the rapidly expanding field of decadal prediction.

Applying the central evaluation system to the three MiKlip hindcast generations has identified a problem with the full-field initializations that to our knowledge has so far escaped attention. While the prototype hindcasts tend to provide the highest skill for North Atlantic Subpolar Gyre SST in later lead years, early lead years display a marked degradation in skill. This degradation is most pronounced in a drop in correlation skill in the initializations with ORAS4 and an

increase in RMSE in the initializations with GECCO2 (Fig. 5). Presumably this skill degradation is related to model drift upon initialization with a state that builds on an incompatible climatology. Figure 5 furthermore illustrates the limitation of our testing procedure with small test ensembles—it is only the full prototype ensemble that identifies the consequences of the drift and forces us to readdress the question of full-field versus anomaly initialization.

As an example of evaluating probabilistic forecasts of discrete events with the RPSS, we analyze whether wind storms related to intense extratropical cyclones occur at a frequency that is either below normal, normal, or above normal for the Northern Hemisphere extended winter season (October through March; Fig. 6; Kruschke et al. 2015). The analysis combines the 29 realizations from all three MiKlip generations available at that time. Using climatology as the reference leads to RPSS-based skill over most of the Northern Hemisphere (not shown; Kruschke et al. 2015). Against the historical simulations as reference, however, additional skill arises in only a few regions, the most prominent of which are the entrance of the North Pacific storm track over eastern Asia and the northwestern Pacific. Similar but less pronounced and less coherent skill enhancement occurs at the entrance of the North Atlantic storm track along the North American east coast and the American sector of the Arctic Ocean (Fig. 6; Kruschke et al. 2015).



**FIG. 5. (a) Correlation and (b) RMSE against HadISST1.1 (Rayner et al. 2003) of the SST index for the North Atlantic Subpolar Gyre (40°–60°N, 0°–60°W) against lead year and for all MiKlip generations and the historical simulations. Baseline0 and baseline1 outperform the historical simulations for almost all lead years. The prototype (pr) hindcasts sometimes provide the highest skill, as is the case for most lead years when using GECCO2 and correlation skill in (a). But sometimes the prototype hindcasts provide the lowest skill, especially for early lead years, as is the case when using ORAS4 and correlation skill in (a) as well as when using either ORAS4 or GECCO2 and the RMSE in (b).**

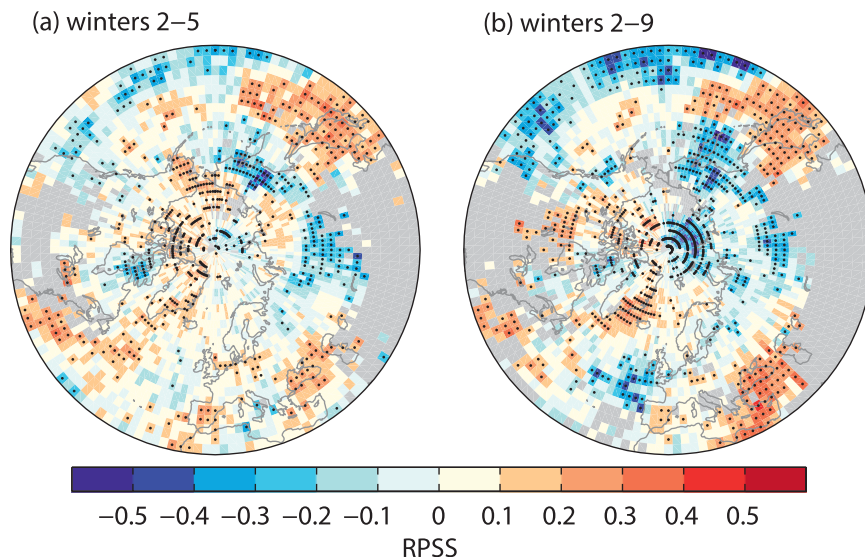


For the analysis shown in Fig. 6, Kruschke et al. (2015) developed and used a bias correction that goes beyond the one recommended in Goddard et al. (2013). The standard correction method is effectively an adjustment of the mean that only depends on lead time. But in a changing climate, model drift following initialization depends also on start year (Kharin et al. 2012). Kruschke et al. (2015) therefore combined the bias correction by Gangstø et al. (2013), which is formulated as a third-order polynomial in lead time, with the drift correction proposed by Kharin et al. (2012) by making the coefficients of the third-order polynomial a linear function of the start year.

We mention here four further examples of evaluating hindcast skill for quantities other than the surface temperature. First, the baseline1-MR version shows prediction skill for the QBO for lead times of up to 4 years. Here, it is essential to use the atmospheric initialization as well as the high vertical resolution in the atmosphere for basic process representation (Pohlmann et al. 2013; see also Scaife et al. 2014b). Second, the MSESS and ensemble reliability have been computed for zonal-mean geopotential height. The only weak dependence of the skill measures on lead time suggests that for geopotential height, changes in external forcing are the main source of skill (Stolzenberger et al. 2015). Third, baseline1 displays significant prediction skill for the AMOC at 26.5°N (Müller et al. 2016, manuscript submitted to *Climate Dyn.*), confirming the earlier results obtained with a system predating the CMIP5 (Matei et al. 2012a), although the physical cause of the prediction skill appears to be different. And fourth, baseline1 shows multiyear potential prediction skill for carbon uptake by the North Atlantic Subpolar Gyre, arising from the improved representation of SST through the initialization (Li et al. 2016).

## PROCESSES AND MODEL DEVELOPMENT.

One MiKlip module aims to understand better the



**FIG. 6.** Hindcast skill for whether wind storms related to intense extratropical cyclones occur at a frequency that is either below normal, normal, or above normal, for the Northern Hemisphere extended winter season (number of tracks within 1,000 km per period Oct–Mar), in a 29-member ensemble constructed from all three MiKlip generations. Skill score is the RPSS; the reference predictions are the historical simulations, and the verification dataset is the ERA reanalyses. (a) Hindcast of winters 2–5 and (b) hindcast of winters 2–9; significant skill scores are indicated ( $\alpha < 5\%$ ) as black dots, and areas of strong inconsistencies between ERA-40 and ERA-Interim are masked out (gray). From Kruschke et al. (2015), reproduced with permission ([www.schweizerbart.de](http://www.schweizerbart.de)).

processes causing decadal variability, to improve existing model components, and to incorporate additional climate subsystems that are relevant for decadal climate predictions. Substantial effort is devoted to exploring the effects of model resolution. For example, a higher-resolution (T106) version of the CMIP3 atmospheric model ECHAM5 revealed that a significant fraction of the convective precipitation over and south of the Gulf Stream can be explained by the variability of the underlying SST, especially in summer (Hand et al. 2014; see also Minobe et al. 2008). Higher horizontal resolution in both atmosphere and ocean is expected to improve the teleconnections between the North Atlantic and Europe (e.g., Minobe et al. 2008; Hand et al. 2014), which are weaker at the T63 atmospheric horizontal resolution used in MiKlip than in reanalyses (e.g., Müller et al. 2012; Ghosh et al. 2016.). Increasing the atmospheric horizontal resolution to T127 is therefore high on MiKlip's list of priorities.

The subpolar North Atlantic and its interaction between gyre and overturning circulations are important for the northward oceanic heat transport and thus for Atlantic warming events such as in the 1990s (Robson et al. 2012a) and the 1920s (Müller

et al. 2015), including their predictions [Robson et al. (2012b) and Müller et al. (2014), respectively]. These results underscore the importance of reducing the misplacement of the Gulf Stream and the North Atlantic Current that is ubiquitous in CMIP5 climate models (e.g., Flato et al. 2013), including the MPI-ESM (Jungclaus et al. 2013).

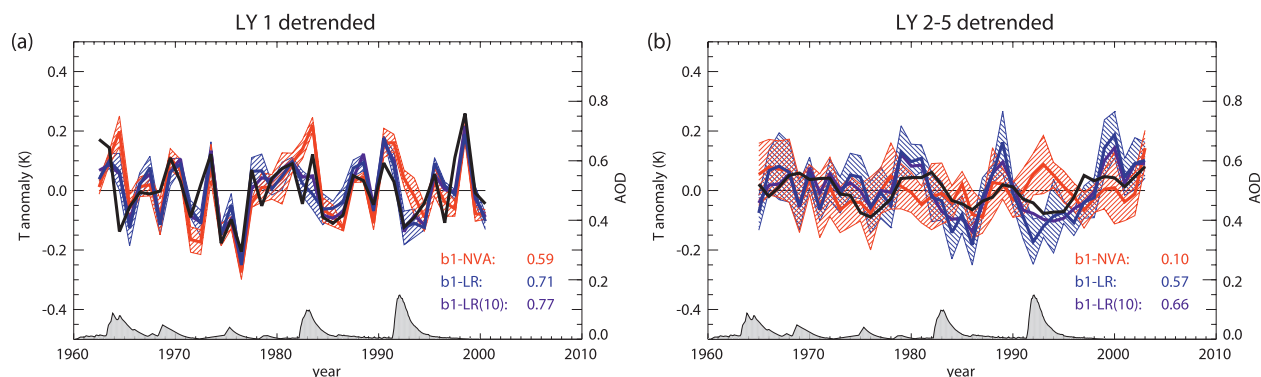
Hindcast skill is markedly degraded by not including the effects of volcanic eruptions (Fig. 7; Timmreck et al. 2016). MiKlip has therefore developed a volcano code package that enables the running of a new ensemble of predictions if a major volcanic eruption occurs in the future. The volcano code package is implemented in a two-step procedure. In the first step, the volcanic radiative forcing is calculated offline with a global aerosol–climate model; in the second step, this forcing is included in the MiKlip system. As a consequence of this two-step procedure, the underlying climate model for producing the predictions remains unchanged, obviating the need to retune the model (Mauritsen et al. 2012) and to create new control and historical simulations.

**DOWNSCALING THE DECADAL PREDICTIONS.** Climate information is often required at a substantially higher spatial resolution than is available from the global climate models, particularly for regional-scale impact studies. The representation of processes such as orographic rain, mesoscale circulations, or wind gusts improves as resolution is refined. For this reason, MiKlip has developed a coordinated regional downscaling component for the

decadal predictions. The two main research questions pursued in MiKlip are (i) whether predictive skill can be found also on the much smaller regional and local scales and (ii) whether the downscaling adds value to the global predictions. The geographical focus lies on Europe and Africa. Because the regional models rely on the global results, there is necessarily some time lag between constructing the global hindcast ensembles and their use in downscaling.

Downscaling implies additional uncertainty (e.g., Räisänen 2007; Flato et al. 2013); therefore, different approaches are employed in MiKlip to assess the robustness of the results. These approaches are coordinated with respect to model grids, initialization, and data processing [analogous to the Coordinated Regional Climate Downscaling Experiment (CORDEX) contribution to CMIP5; e.g., Kotlarski et al. 2014]. For Europe, the ensemble consists of the two regional climate models (RCMs) Consortium for Small-Scale Modelling in Climate Mode (COSMO-CLM or CCLM; Rockel et al. 2008) and Regional-Scale Model (REMO; Jacob 2001), and a statistical–dynamical method. For Africa, three RCMs are used: CCLM, REMO, and Weather Research and Forecasting (WRF) Model (Skamarock and Klemp 2008).

The regionalization for Europe maintains or slightly enhances the skill inherited from the baseline global hindcasts for annual-mean surface temperature (Fig. 8). Given the user orientation of downscaled predictions, we show here the combined skill from forcing changes and initialized internal variability; skill score is MSESS evaluated against the European daily high-



**FIG. 7.** Detrended time series of hindcast and observed globally averaged surface temperature from 1962 to 2004; the figure shows anomalies relative to the mean over this period for HadCRUT3v observations (Brohan et al. 2006; black), baseline-LR (blue), and baseline-LR without volcanic eruptions (b1-NVA; red). The blue and red curves are each based on three realizations. The purple line indicates the 10-ensemble-member mean of baseline-LR. The standard deviation is indicated by the hatched areas. (a) Lead year 1 and (b) lead years 2–5. The numbers indicate the anomaly correlation coefficient between the hindcasts and the observations over the whole period. The gray shaded region right above the x axis shows the time series of annually averaged stratospheric aerosol optical depth (Stenchikov et al. 1998; and updates). From Timmreck et al. (2016), reproduced with permission.

resolution gridded dataset (E-OBS; Haylock et al. 2008), with climatology as the reference forecast. The RCM ensemble consists of simulations with CCLM as well as with REMO, and it maintains the skill in western and southern Europe and shows an increase in parts of central, eastern, and northern Europe (Fig. 8).

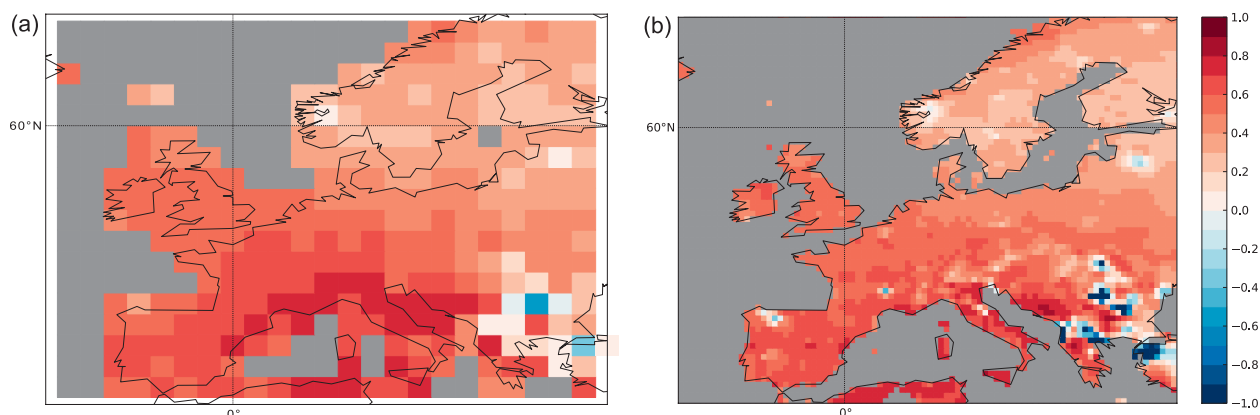
Added value of the downscaling has been found for strong precipitation events over central Europe; the RCM CCLM clearly outperforms the baseline0 global model in the representation of the frequency of days with precipitation larger than about 20 mm day<sup>-1</sup> (not shown; Mieruch et al. 2014). Furthermore, while the global model ensemble is overconfident (ensemble spread smaller than the error, a feature that is ever more pronounced with increasing precipitation intensity), the regional model ensemble is reliable out to very large intensities.

A statistical–dynamical downscaling approach comprising a combination of weather typing and CCLM simulations has been used to explore the predictability of wind energy output over central Europe (Reyers et al. 2015). The skill score used is the MESS, the reference prediction is the downscaled historical simulation, and the verification dataset is the downscaled wind energy output of ERA-Interim for the period 1979–2010. While no skill is found for baseline0, positive skill is obtained for short forecast periods of baseline1 and prototype, particularly over central Europe; prototype GECCO2 outperforms all other systems over Poland for lead years 2–5 (Fig. 9). Hindcast skill is highest for autumn and lowest for summer over central Europe (not shown), indicating a clear dependency of the predictive skill on season (Moemken et al. 2016).

**DISCUSSION AND CONCLUSIONS.** MiKlip is well poised to deliver its decadal prediction and evaluation systems to the DWD for operational use by 2019. Placing a single global prediction system in the focus of a major research effort has demonstrated benefits such as the rapid development of alternative initialization strategies, sophisticated evaluation methods for quantities beyond the surface temperature, and regional applications of the global predictions. Such rapid progress would have been impossible at any single institution in Germany, no matter how scientifically powerful or well-funded.

At least five major issues remain unsettled and must be tackled by MiKlip in the coming years:

- 1) We have not yet converged on a best initialization procedure of our prediction ensemble. Some hindcasts suffer from degraded skill right after initialization, in particular when full-field initialization is used. This effect presumably is related to using an assimilation model, either statistical or dynamical, that is different from the model used in the hindcasts (Kröger et al. 2012). Furthermore, it is unsatisfactory that our initial condition ensemble is unable to explore the full uncertainty range of the initial ocean state.
- 2) The teleconnections between SST and surface temperature over land are not robust enough in our model. While MiKlip has successfully reproduced the observed connection between the SST in the tropical Atlantic and the West African monsoon (Paeth et al. 2016, manuscript submitted to *Meteorologische Zeitschrift*), prediction



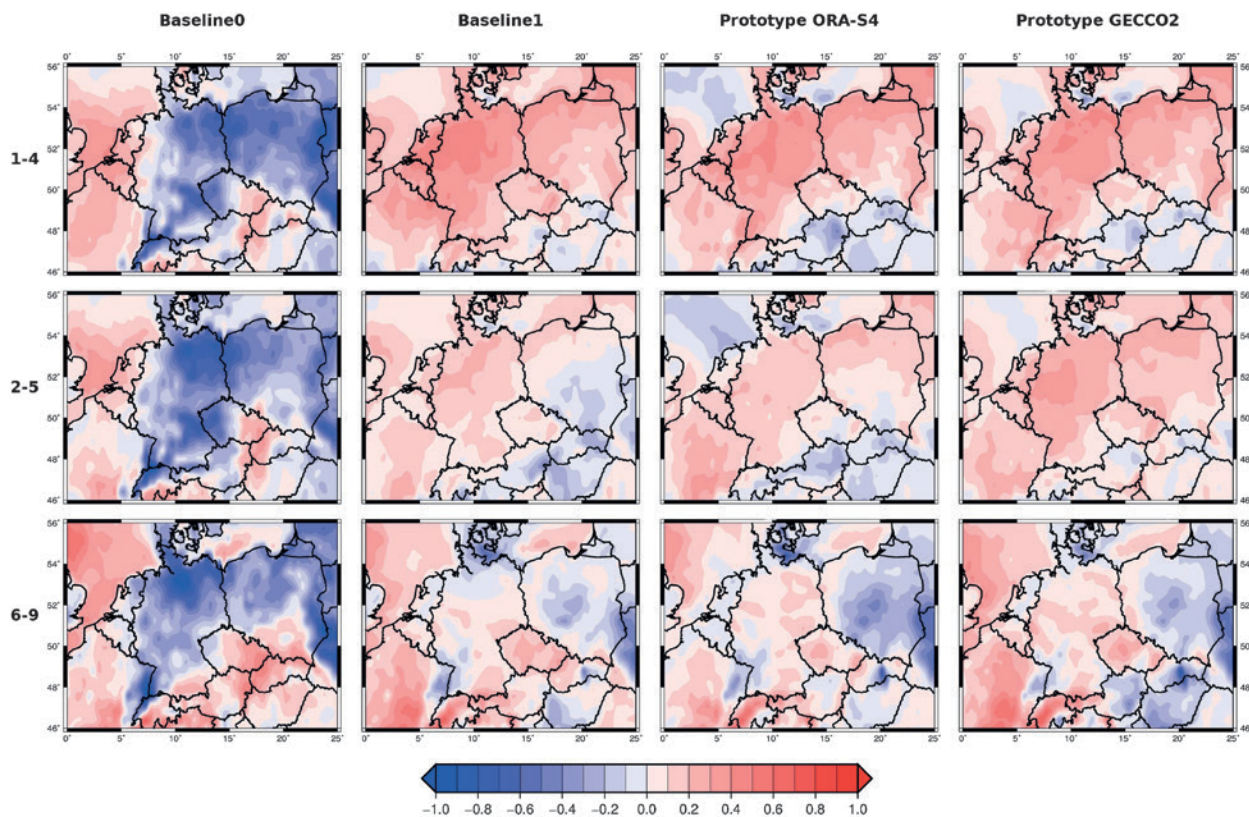
**FIG. 8.** Ensemble-mean hindcast skill over Europe for near-surface air temperature for lead time 2–5 yr, starting yearly from 1961 to 2004. Skill score is MESS evaluated against E-OBS (Haylock et al. 2008), with climatology as the reference forecast; positive values indicate better skill than climatology. (a) MPI-ESM-LR baseline1; (b) RCM ensemble (combined from CCLM and REMO). Given the user orientation of downscaled predictions, we show here the combined skill from forcing changes and initialized internal variability.

skill for North Atlantic SST translates into only some, but not sufficient, skill over Europe (Müller et al. 2012). The required higher-resolution version of MPI-ESM has until recently not been available, owing to some unrealistic features in an earlier control run (J. Jungclaus 2014, personal communication). These problems have now been overcome, and we will perform the next set of production runs with an atmospheric model with resolution T127 (MPI-ESM-HR).

- 3) The availability of the MPI-ESM-HR brings into even sharper relief the computing resource issue that we already faced when applying the MR version of our system. Because higher resolution usually implies smaller possible ensemble size, we experience a palpable trade-off between more realistic representation of physical processes on the one hand and the translation of this representation into prediction skill on the other hand. With a new computer available to MiKlip since July 2015, the competition for resources between

resolution and ensemble size has subsided somewhat, but in the foreseeable future hindcasts with MPI-ESM-HR will be limited to an ensemble size of 10.

- 4) When starting MiKlip, we underestimated the difficulty of implementing suggested model improvements. Any modification to the climate model itself requires a retuning (e.g., Mauritsen et al. 2012), a new control run with constant forcing to make sure the model simulates a stable climate, and a new ensemble of historical runs as a reference for assessing skill enhancement through initialization. Being tied to the general MPI-ESM development implies that the cycle of model versions rests outside of MiKlip's immediate control and occurs in intervals longer than sometimes desired by MiKlip. On the other hand, MiKlip does not command the personnel resources needed to maintain an independent climate model, and even if it did, separating its model development from that of the MPI-ESM would not use resources efficiently—



**FIG. 9.** MSESS for wind energy output for (top) years 1–4, (middle) years 2–5, and (bottom) years 6–9 of the (left column) baseline0 3-member ensemble mean, (second column) baseline1 10-member ensemble mean, (third column) prototype ORAS4 10-member subensemble mean, and (right column) prototype GECCO2 10-member subensemble mean. Reference prediction is the ensemble mean of the uninitialized historical runs, using 3 realizations for baseline0 and 10 members for the other generations; verification dataset is the wind energy output downscaled from ERA-Interim (Dee et al. 2011). Adapted from Moemken et al. (2016).

MiKlip would maintain a full-blown climate model for decadal prediction alone.

For generational cycles of the prediction system that are defined not through different model versions but through different initialization procedures, a much faster turnover can be implemented. The 18-month turnover originally envisioned in MiKlip, however, proved to be overambitious for a sustained mode of operation. We therefore decided not to produce a set of hindcasts during development stage 3 and have instead focused our effort on a comprehensive evaluation of the prototype system. A sustained 18-month turnover would imply that we could never explore the full implication of a generation of hindcasts, including the effects on downscaling, before designing the generation after. We thus tentatively recommend for later operational use to allow for a more relaxed cycle of prediction system generations, with intervals of 2–3 years rather than 18 months.

- 5) We have so far focused almost exclusively on evaluating the hindcasts and not on constructing and issuing our own exploratory forecasts, although we do participate in the multimodel real-time decadal prediction exercise led by the Hadley Centre (Smith et al. 2013b). We have also started a dialogue with potential users of the MiKlip forecasts and have now added subprojects that develop such a dialogue systematically. Issuing our own forecasts requires further exploration of how to communicate the strengths and weaknesses of the forecast in a manner both accurate and easy to grasp. MiKlip plans to tackle this challenge over the coming years because without this communication component an operational system would remain incomplete.

**ACKNOWLEDGMENTS.** MiKlip is funded by the German Federal Ministry for Education and Research (BMBF) under grant agreements 01LP11nnx, where nn ranges from 04 to 70 and x ranges from A to F. All simulations were carried out at the German Climate Computing Centre (DKRZ), which also provided all major data services. We thank Bjorn Stevens, the anonymous reviewers, and Editor Michael Alexander for comments on an earlier version of the manuscript.

## REFERENCES

- Andersson, A., M. Ziese, F. Dietzsch, M. Schröder, A. Becker, and K. Schamm, 2016a: HOAPS/GPCC global daily precipitation data record with uncertainty estimates using satellite and gauge based observations at 1.0°. Deutscher Wetterdienst, accessed 1 June 2016, doi:10.5676/DWD\_CDC/HOGP\_100/V001.
- , —, —, —, —, and —, 2016b: HOAPS/GPCC global daily precipitation data record with uncertainty estimates using satellite and gauge based observations at 2.5°. Deutscher Wetterdienst, accessed 1 June 2016, doi:10.5676/DWD\_CDC/HOGP\_250/V001.
- Balmaseda, M. A., K. Mogensen, and A. T. Weaver, 2013: Evaluation of the ECMWF ocean reanalysis system ORAS4. *Quart. J. Roy. Meteor. Soc.*, **139**, 1132–1161, doi:10.1002/qj.2063.
- Bellucci, A., and Coauthors, 2015a: Advancements in decadal climate predictability: The role of nonoceanic drivers. *Rev. Geophys.*, **53**, 165–202, doi:10.1002/2014RG000473.
- , and Coauthors, 2015b: An assessment of a multimodel ensemble of decadal climate predictions. *Climate Dyn.*, **44**, 2787–2806, doi:10.1007/s00382-014-2164-y.
- Brohan, P., J. J. Kennedy, I. Harris, S. F. B. Tett, and P. D. Jones, 2006: Uncertainty estimates in regional and global observed temperature changes: A new data set from 1850. *J. Geophys. Res.*, **111**, D12106, doi:10.1029/2005JD006548.
- Brune, S., L. Nerger, and J. Baehr, 2015: Assimilation of oceanic observations in a global coupled Earth system model with the SEIK filter. *Ocean Modell.*, **96**, 254–264, doi:10.1016/j.ocemod.2015.09.011.
- Cunningham, S. A., and Coauthors, 2007: Temporal variability of the Atlantic meridional overturning circulation at 26.5°N. *Science*, **317**, 935–938, doi:10.1126/science.1141304.
- Dee, D. P., and Coauthors, 2011: The ERA-Interim reanalysis: Configuration and performance of the data assimilation system. *Quart. J. Roy. Meteor. Soc.*, **137**, 553–597, doi:10.1002/qj.828.
- Du, H., F. J. Doblas-Reyes, J. García-Serrano, V. Guemas, Y. Soufflet, and B. Wouters, 2012: Sensitivity of decadal predictions to the initial atmospheric and oceanic perturbations. *Climate Dyn.*, **39**, 2013–2023, doi:10.1007/s00382-011-1285-9.
- Flato, G., and Coauthors, 2013: Evaluation of climate models. *Climate Change 2013: The Physical Science Basis*, T. F. Stocker et al., Eds., Cambridge University Press, 741–866.
- Gangstø, R., A. P. Weigel, M. A. Liniger, and C. Appenzeller, 2013: Methodological aspects of the validation of decadal predictions. *Climate Res.*, **55**, 181–200, doi:10.3354/cr01135.
- Ghosh, R., W. A. Müller, J. Bader, and J. Baehr, 2016: Impact of observed North Atlantic multidecadal variations to European summer climate: A quasi-

- geostrophic pathway. *Climate Dyn.*, doi:10.1007/s00382-016-3283-4.
- Giorgetta, M. A., and Coauthors, 2013: 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.*, **5**, 572–597, doi:10.1002/jame.20038.
- Goddard, L., and Coauthors, 2013: A verification framework for interannual-to-decadal predictions experiments. *Climate Dyn.*, **40**, 245–272, doi:10.1007/s00382-012-1481-2.
- Guemas, V., F. J. Doblas-Reyes, F. Lienert, Y. Soufflet, and H. Du, 2012: Identifying the causes of the poor decadal climate prediction skill over the North Pacific. *J. Geophys. Res.*, **117**, D20111, doi:10.1029/2012JD018004.
- Hand, R., N. Keenlyside, N.-E. Omrani, and M. Latif, 2014: Simulated response to inter-annual SST variations in the Gulf Stream region. *Climate Dyn.*, **42**, 715–731, doi:10.1007/s00382-013-1715-y.
- Harris, I., P. D. Jones, T. J. Osborn, and D. H. Lister, 2014: Updated high-resolution grids of monthly climatic observations—The CRU TS3.10 dataset. *Int. J. Climatol.*, **34**, 623–642, doi:10.1002/joc.3711.
- Haylock, M. R., N. Hofstra, A. M. G. K. Tank, E. J. Klok, P. D. Jones, and M. New, 2008: A European daily high-resolution gridded data set of surface temperature and precipitation for 1950–2006. *J. Geophys. Res.*, **113**, D20119, doi:10.1029/2008JD010201.
- Hermanson, L., R. Eade, N. H. Robinson, N. J. Dunstone, M. B. Andrews, J. R. Knight, A. A. Scaife, and D. M. Smith, 2014: Forecast cooling of the Atlantic Subpolar Gyre and associated impacts. *Geophys. Res. Lett.*, **41**, 5167–5174, doi:10.1002/2014GL060420.
- Illing, S., C. Kadow, O. Kunst, and U. Cubasch, 2014: MurCSS: A tool for standardized evaluation of decadal hindcast systems. *J. Open Res. Software*, **2**, e24, doi:10.5334/jors.bf.
- Ishii, M., and M. Kimoto, 2009: Reevaluation of historical ocean heat content variations with time-varying XBT and MBT depth bias corrections. *J. Oceanogr.*, **65**, 287–299, doi:10.1007/s10872-009-0027-7.
- Jacob, D., 2001: A note to the simulation of the annual and inter-annual variability of the water budget over the Baltic Sea drainage basin. *Meteor. Atmos. Phys.*, **77**, 61–73, doi:10.1007/s007030170017.
- Jones, P. D., D. H. Lister, T. J. Osborn, C. Harpham, M. Salmon, and C. P. Morice, 2012: Hemispheric and large-scale land-surface air temperature variations: An extensive revision and an update to 2010. *J. Geophys. Res.*, **117**, D05127, doi:10.1029/2011JD017139.
- Jungclaus, J. H., and Coauthors, 2013: Characteristics of the ocean simulations in the Max Planck Institute Ocean Model (MPIOM) the ocean component of the MPI-Earth system model. *J. Adv. Model. Earth Syst.*, **5**, 422–446, doi:10.1002/jame.20023.
- Kadow, C., S. Illing, O. Kunst, H. W. Rust, H. Pohlmann, W. A. Müller, and U. Cubasch, 2015: Evaluation of forecasts by accuracy and spread in the MiKlip decadal climate prediction system. *Meteor. Z.*, doi:10.1127/metz/2015/0639.
- Kalnay, E., and Coauthors, 1996: The NCEP/NCAR 40-Year Reanalysis Project. *Bull. Amer. Meteor. Soc.*, **77**, 437–471, doi:10.1175/1520-0477(1996)077<0437: TNYRP>2.0.CO;2.
- Keenlyside, N. S., M. Latif, J. Jungclaus, L. Kornblueh, and E. Roeckner, 2008: Advancing decadal-scale climate prediction in the North Atlantic sector. *Nature*, **453**, 84–88, doi:10.1038/nature06921.
- Kharin, V. V., G. J. Boer, W. J. Merryfield, J. F. Scinocca, and W. S. Lee, 2012: Statistical adjustment of decadal predictions in a changing climate. *Geophys. Res. Lett.*, **39**, L19705, doi:10.1029/2012GL052664.
- Kirtman, B., and Coauthors, 2013: Near-term climate change: Projections and predictability. *Climate Change 2013: The Physical Science Basis*, T. F. Stocker et al., Eds., Cambridge University Press, 953–1028.
- Köhl, A., 2015: Evaluation of the GECCO2 ocean synthesis: Transports of volume, heat and freshwater in the Atlantic. *Quart. J. Roy. Meteor. Soc.*, **141**, 166–181, doi:10.1002/qj.2347.
- Kotlarski, S., and Coauthors, 2014: Regional climate modeling on European scales: A joint standard evaluation of the EURO-CORDEX RCM ensemble. *Geosci. Model Dev.*, **7**, 1297–1333, doi:10.5194/gmd-7-1297-2014.
- Kröger, J., W. A. Müller, and J. S. von Storch, 2012: Impact of different ocean reanalyses on decadal climate prediction. *Climate Dyn.*, **39**, 795–810, doi:10.1007/s00382-012-1310-7.
- Kruschke, T., H. W. Rust, C. Kadow, W. A. Müller, H. Pohlmann, G. C. Leckebusch, and U. Ulbrich, 2015: Probabilistic evaluation of decadal prediction skill regarding Northern Hemisphere winter storms. *Meteor. Z.*, doi:10.1127/metz/2015/0641.
- Kumar, A., A. G. Barnston, and M. P. Hoerling, 2001: Seasonal predictions, probabilistic verifications, and ensemble size. *J. Climate*, **14**, 1671–1676, doi:10.1175/1520-0442(2001)014<1671:SPPVAE>2.0.CO;2.
- Lee, T., D. E. Waliser, J.-L. F. Li, F. W. Landerer, and M. M. Gierach, 2013: Evaluation of CMIP3 and CMIP5 wind stress climatology using satellite measurements and atmospheric reanalysis products. *J. Climate*, **26**, 5810–5826, doi:10.1175/JCLI-D-12-00591.1.
- Levitus, S., J. I. Antonov, T. P. Boyer, R. A. Locarnini, H. E. Garcia, and A. V. Mishonov, 2009: Global ocean

- heat content 1955–2008 in light of recently revealed instrumentation problems. *Geophys. Res. Lett.*, **36**, L07608, doi:10.1029/2008GL037155.
- Li, H., T. Ilyina, W. A. Müller, and F. Sienz, 2016: Decadal predictions of the North Atlantic CO<sub>2</sub> uptake. *Nat. Commun.*, **7**, 11076, doi:10.1038/ncomms11076.
- Marini, C., I. Polkova, A. Köhl, and D. Stammer, 2016: A comparison of two ensemble generation methods using oceanic singular vectors and atmospheric lagged initialization for decadal climate prediction. *Mon. Wea. Rev.*, **144**, 2719–2738, doi:10.1175/MWR-D-15-0350.1.
- Matei, D., J. Baehr, J. H. Jungclaus, H. Haak, W. A. Müller, and J. Marotzke, 2012a: Multiyear prediction of monthly mean Atlantic meridional overturning circulation at 26.5°N. *Science*, **335**, 76–79, doi:10.1126/science.1210299.
- , H. Pohlmann, J. Jungclaus, W. Müller, H. Haak, and J. Marotzke, 2012b: Two tales of initializing decadal climate prediction experiments with the ECHAM5/MPI-OM model. *J. Climate*, **25**, 8502–8523, doi:10.1175/JCLI-D-11-00633.1.
- Mauritsen, T., and Coauthors, 2012: Tuning the climate of a global model. *J. Adv. Model. Earth Syst.*, **4**, M00A01, doi:10.1029/2012MS000154.
- McCarthy, G., and Coauthors, 2012: Observed interannual variability of the Atlantic meridional overturning circulation at 26.5°N. *Geophys. Res. Lett.*, **39**, L19609, doi:10.1029/2012GL052933.
- McGregor, S., A. Sen Gupta, and M. H. England, 2012: Constraining wind stress products with sea surface height observations and implications for Pacific Ocean sea level trend attribution. *J. Climate*, **25**, 8164–8176, doi:10.1175/JCLI-D-12-00105.1.
- Meehl, G. A., and Coauthors, 2014: Decadal climate prediction: An update from the trenches. *Bull. Amer. Meteor. Soc.*, **95**, 243–267, doi:10.1175/BAMS-D-12-00241.1.
- Mieruch, S., H. Feldmann, G. Schädler, C. J. Lenz, S. Kothe, and C. Kottmeier, 2014: The regional MiKlip decadal forecast ensemble for Europe: The added value of downscaling. *Geosci. Model Dev.*, **7**, 2983–2999, doi:10.5194/gmd-7-2983-2014.
- Minobe, S., A. Kuwano-Yoshida, N. Komori, S. P. Xie, and R. J. Small, 2008: Influence of the Gulf Stream on the troposphere. *Nature*, **452**, 206–209, doi:10.1038/nature06690.
- Mochizuki, T., and Coauthors, 2010: Pacific decadal oscillation hindcasts relevant to near-term climate prediction. *Proc. Natl. Acad. Sci. USA*, **107**, 1833–1837, doi:10.1073/pnas.0906531107.
- Moemken, J., M. Reyers, B. Buldmann, and J. G. Pinto, 2016: Decadal predictability of regional scale wind speed and wind energy potentials over central Europe. *Tellus*, **68A**, 29199, doi:10.3402/tellusa.v68.29199.
- Molteni, F., R. Buizza, T. N. Palmer, and T. Petroliagis, 1996: The ECMWF ensemble prediction system: Methodology and validation. *Quart. J. Roy. Meteor. Soc.*, **122**, 73–119, doi:10.1002/qj.49712252905.
- Müller, W. A., and Coauthors, 2012: Forecast skill of multi-year seasonal means in the decadal prediction system of the Max Planck Institute for Meteorology. *Geophys. Res. Lett.*, **39**, L22707, doi:10.1029/2012GL053326.
- , H. Pohlmann, F. Sienz, and D. Smith, 2014: Decadal climate predictions for the period 1901–2010 with a coupled climate model. *Geophys. Res. Lett.*, **41**, 2100–2107, doi:10.1002/2014GL059259.
- , and Coauthors, 2015: A twentieth-century reanalysis forced ocean model to reconstruct the North Atlantic climate variation during the 1920s. *Climate Dyn.*, **44**, 1935–1955, doi:10.1007/s00382-014-2267-5.
- Murphy, A. H., 1988: Skill scores based on the mean-square error and their relationships to the correlation coefficient. *Mon. Wea. Rev.*, **116**, 2417–2425, doi:10.1175/1520-0493(1988)116<2417:SSBOTM>2.CO;2.
- Murphy, J. M., 1990: Assessment of the practical utility of extended range ensemble forecasts. *Quart. J. Roy. Meteor. Soc.*, **116**, 89–125, doi:10.1002/qj.49711649105.
- Pham, D. T., J. Verron, and M. C. Roubaud, 1998: A singular evolutive extended Kalman filter for data assimilation in oceanography. *J. Mar. Syst.*, **16**, 323–340, doi:10.1016/S0924-7963(97)00109-7.
- Pohlmann, H., J. H. Jungclaus, A. Köhl, D. Stammer, and J. Marotzke, 2009: Initializing decadal climate predictions with the GECCO oceanic synthesis: Effects on the North Atlantic. *J. Climate*, **22**, 3926–3938, doi:10.1175/2009JCLI2535.1.
- , and Coauthors, 2013: Improved forecast skill in the tropics in the new MiKlip decadal climate predictions. *Geophys. Res. Lett.*, **40**, 5798–5802, doi:10.1002/2013GL058051.
- Polkova, I., A. Köhl, and D. Stammer, 2014: Impact of initialization procedures on the predictive skill of a coupled ocean–atmosphere model. *Climate Dyn.*, **42**, 3151–3169, doi:10.1007/s00382-013-1969-4.
- Räsänen, J., 2007: How reliable are climate models? *Tellus*, **59A**, 2–29, doi:10.1111/j.1600-0870.2006.00211.x.
- Rayner, N. A., D. E. Parker, E. B. Horton, C. K. Folland, L. V. Alexander, D. P. Rowell, E. C. Kent, and A. Kaplan, 2003: Global analyses of sea surface

- temperature, sea ice, and night marine air temperature since the late nineteenth century. *J. Geophys. Res.*, **108**, 4407, doi:10.1029/2002JD002670.
- Reyers, M., J. G. Pinto, and J. Moemken, 2015: Statistical-dynamical downscaling for wind energy potentials: Evaluation and applications to decadal hindcasts and climate change projections. *Int. J. Climatol.*, **35**, 229–244, doi:10.1002/joc.3975.
- Robson, J., R. Sutton, K. Lohmann, D. Smith, and M. D. Palmer, 2012a: Causes of the rapid warming of the North Atlantic Ocean in the mid-1990s. *J. Climate*, **25**, 4116–4134, doi:10.1175/JCLI-D-11-00443.1.
- , —, and D. M. Smith, 2012b: Initialized decadal predictions of the rapid warming of the North Atlantic Ocean in the mid 1990s. *Geophys. Res. Lett.*, **39**, L19713, doi:10.1029/2012GL053370.
- Rockel, B., A. Will, and A. Hense, 2008: The regional climate model COSMO-CLM (CCLM). *Meteor. Z.*, **17**, 347–348, doi:10.1127/0941-2948/2008/0309.
- Romanova, V., and A. Hense, 2015: Anomaly transform methods based on total energy and ocean heat content norms for generating ocean dynamic disturbances for ensemble climate forecasts. *Climate Dyn.*, doi:10.1007/s00382-015-2567-4, in press.
- Scaife, A. A., and Coauthors, 2014a: Skillful long-range prediction of European and North American winters. *Geophys. Res. Lett.*, **41**, 2514–2519, doi:10.1002/2014GL059637.
- , and Coauthors, 2014b: Predictability of the quasi-biennial oscillation and its northern winter teleconnection on seasonal to decadal timescales. *Geophys. Res. Lett.*, **41**, 1752–1758, doi:10.1002/2013GL059160.
- Schamm, K., M. Ziese, A. Becker, P. Finger, A. Meyer-Christoffer, U. Schneider, M. Schröder, and P. Stender, 2014: Global gridded precipitation over land: A description of the new GPCC first guess daily product. *Earth Syst. Sci. Data*, **6**, 49–60, doi:10.5194/essd-6-49-2014.
- Sienz, F., H. Pohlmann, and W. A. Müller, 2016: Ensemble size impact on the decadal predictive skill assessment. *Meteor. Z.*, doi:10.1127/metz/2016/0670.
- Skamarock, W. C., and J. B. Klemp, 2008: A time-split nonhydrostatic atmospheric model for weather research and forecasting applications. *J. Comput. Phys.*, **227**, 3465–3485, doi:10.1016/j.jcp.2007.01.037.
- Smith, D. M., S. Cusack, A. W. Colman, C. K. Folland, G. R. Harris, and J. M. Murphy, 2007: Improved surface temperature prediction for the coming decade from a global climate model. *Science*, **317**, 796–799, doi:10.1126/science.1139540.
- , R. Eade, and H. Pohlmann, 2013a: A comparison of full-field and anomaly initialization for seasonal to decadal climate prediction. *Climate Dyn.*, **41**, 3325–3338, doi:10.1007/s00382-013-1683-2.
- , and Coauthors, 2013b: Real-time multi-model decadal climate predictions. *Climate Dyn.*, **41**, 2875–2888, doi:10.1007/s00382-012-1600-0.
- Stenchikov, G. L., I. Kirchner, A. Robock, H. F. Graf, J. C. Antuna, R. G. Grainger, A. Lambert, and L. Thomason, 1998: Radiative forcing from the 1991 Mount Pinatubo volcanic eruption. *J. Geophys. Res.*, **103**, 132837–132857, doi:10.1029/98JD00693.
- Stolzenberger, S., R. Glowienka-Hense, T. Spanghel, M. Schröder, A. Mazurkiewicz, and A. Hense, 2015: Revealing skill of the MiKliP decadal prediction systems by three-dimensional probabilistic evaluation. *Meteor. Z.*, doi:10.1127/metz/2015/0606.
- Taylor, K. E., R. J. Stouffer, and G. A. Meehl, 2012: An overview of CMIP5 and the experiment design. *Bull. Amer. Meteor. Soc.*, **93**, 485–498, doi:10.1175/BAMS-D-11-00094.1.
- Thoma, M., R. J. Greatbatch, C. Kadow, and R. Gerdes, 2015: Decadal hindcasts initialized using observed surface wind stress: Evaluation and prediction out to 2024. *Geophys. Res. Lett.*, **42**, 6454–6461, doi:10.1002/2015GL064833.
- Timmreck, C., H. Pohlmann, C. Kadow, and S. Illing, 2016: The impact of stratospheric volcanic aerosol on decadal scale predictability. *Geophys. Res. Lett.*, **43**, 834–842, doi:10.1002/2015GL067431.
- Uppala, S. M., and Coauthors, 2005: The ERA-40 reanalysis. *Quart. J. Roy. Meteor. Soc.*, **131**, 2961–3012, doi:10.1256/qj.04.176.
- Wei, M. Z., Z. Toth, R. Wobus, Y. J. Zhu, C. H. Bishop, and X. G. Wang, 2006: Ensemble transform Kalman filter-based ensemble perturbations in an operational global prediction system at NCEP. *Tellus*, **58A**, 28–44, doi:10.1111/j.1600-0870.2006.00159.x.
- WMO, 2011: Climate knowledge for action: A global framework for climate services—Empowering the most vulnerable. WMO Rep. 1065, 247 pp.
- Zhang, L., H. Dobslaw, C. Dahle, I. Sasgen, and M. Thomas, 2015: Validation of MPI-ESM decadal hindcast experiments with terrestrial water storage variations as observed by the GRACE satellite mission. *Meteor. Z.*, doi:10.1127/metz/2015/0596.