

# *North Atlantic climate far more predictable than models imply*

Article

Accepted Version

Smith, D. M., Scaife, A. A., Eade, R., Athanasiadis, P., Bellucci, A., Bethke, I., Bilbao, R., Borchert, L. F., Caron, L.-P., Counillon, F., Danabasoglu, G., Delworth, T., Doblas-Reyes, F. J., Dunstone, N. J., Estella-Perez, V., Flavoni, S., Hermanson, L., Keenlyside, N., Kharin, V., Kimoto, M., Merryfield, W. J., Mignot, J., Mochizuki, T., Modali, K., Monerie, P.-A. ORCID: <https://orcid.org/0000-0002-5304-9559>, Müller, W. A., Nicoli, D., Ortega, P., Pankatz, K., Pohlmann, H., Robson, J. ORCID: <https://orcid.org/0000-0002-3467-018X>, Ruggieri, P., Sospedra-Alfonso, R., Swingedouw, D., Wang, Y., Wild, S., Yeager, S., Yang, X. and Zhang, L. (2020) North Atlantic climate far more predictable than models imply. *Nature*, 583. 7818. ISSN 0028-0836 doi: <https://doi.org/10.1038/s41586-020-2525-0> Available at <https://centaur.reading.ac.uk/91527/>

It is advisable to refer to the publisher's version if you intend to cite from the work. See [Guidance on citing](#).

To link to this article DOI: <http://dx.doi.org/10.1038/s41586-020-2525-0>

Publisher: Nature Publishing Group

including copyright law. Copyright and IPR is retained by the creators or other copyright holders. Terms and conditions for use of this material are defined in the [End User Agreement](#).

[www.reading.ac.uk/centaur](http://www.reading.ac.uk/centaur)

## **CentAUR**

Central Archive at the University of Reading

Reading's research outputs online

# Climate models underpredict North Atlantic atmospheric circulation changes

D. M. Smith<sup>1</sup>, A. A. Scaife<sup>1,2</sup>, R. Eade<sup>1</sup>, P. Athanasiadis<sup>3</sup>, A. Bellucci<sup>3</sup>, I. Bethke<sup>4,5</sup>, R. Bilbao<sup>6</sup>, L. F. Borchert<sup>7</sup>, L.-P. Caron<sup>6</sup>, F. Counillon<sup>4,5</sup>, G. Danabasoglu<sup>8</sup>, T. Delworth<sup>9</sup>, F. J. Doblas-Reyes<sup>6,10</sup>, N. J. Dunstone<sup>1</sup>, V. Estella-Perez<sup>7</sup>, S. Flavoni<sup>7</sup>, L. Hermanson<sup>1</sup>, N. Keenlyside<sup>4,5</sup>, V. Kharin<sup>11</sup>, M. Kimoto<sup>12</sup>, W. J. Merryfield<sup>11</sup>, J. Mignot<sup>7</sup>, T. Mochizuki<sup>13,14</sup>, K. Modali<sup>15</sup>, P.-A. Monerie<sup>16</sup>, W. A. Müller<sup>15</sup>, D. Nicolí<sup>3</sup>, P. Ortega<sup>6</sup>, K. Pankatz<sup>17</sup>, H. Pohlmann<sup>15,17</sup>, J. Robson<sup>16</sup>, P. Ruggieri<sup>3</sup>, R. Sospedra-Alfonso<sup>11</sup>, D. Swingedouw<sup>18</sup>, Y. Wang<sup>4</sup>, S. Wild<sup>6</sup>, S. Yeager<sup>8</sup>, X. Yang<sup>9</sup> and L. Zhang<sup>9</sup>

<sup>1</sup>*Met Office Hadley Centre, FitzRoy Road, Exeter EX1 3PB, UK*

<sup>2</sup>*College of Engineering, Mathematics and Physical Sciences, Exeter University, UK*

<sup>3</sup>*Centro Euro-Mediterraneo sui Cambiamenti Climatici, Bologna, Italy*

<sup>4</sup>*Nansen Environmental and Remote Sensing Center and Bjerknes Centre for Climate Research, Bergen, Norway*

<sup>5</sup>*Geophysical Institute, University of Bergen and Bjerknes Centre for Climate Research, Bergen, Norway*

<sup>6</sup>*Barcelona Supercomputing Center, Jordi Girona 29 - 08034 Barcelona, Spain*

<sup>7</sup>*Sorbonne Universités, LOCEAN Laboratory, Institut Pierre Simon Laplace (IPSL), Paris, France*

<sup>8</sup>*National Center for Atmospheric Research, Boulder, CO, USA*

<sup>9</sup>*Geophysical Fluid Dynamics Laboratory, Princeton University, Princeton, NJ, USA*

<sup>10</sup>*ICREA, Barcelona, Spain*

<sup>11</sup>*Canadian Centre for Climate Modelling and Analysis, Environment and Climate Change Canada,*

22 *Victoria, British Columbia, Canada*

23 <sup>12</sup>*Atmosphere and Ocean Research Institute, University of Tokyo, Kashiwa, Japan*

24 <sup>13</sup>*Department of Earth and Planetary Sciences, Kyushu University, Fukuoka, Japan*

25 <sup>14</sup>*Japan Agency for Marine-Earth Science and Technology, Yokohama, Japan*

26 <sup>15</sup>*Max-Planck-Institut für Meteorologie, Bundesstraße 53, 20146 Hamburg, Germany*

27 <sup>16</sup>*National Centre for Atmospheric Science, Department of Meteorology, University of Reading,*  
28 *Reading RG6 6BB, UK*

29 <sup>17</sup>*Deutscher Wetterdienst, Bernhard-Nocht-Str. 76, Hamburg, Germany*

30 <sup>18</sup>*CNRS-EPOC, Université de Bordeaux, Pessac, France*

31 *Corresponding author: Doug Smith, [doug.smith@metoffice.gov.uk](mailto:doug.smith@metoffice.gov.uk)*

## 32 **Abstract**

33 **Quantifying signals and uncertainties in climate models is essential for climate change de-**  
34 **tection, attribution, prediction and projection<sup>1-3</sup>. Although inter-model agreement is high**  
35 **for large-scale temperature signals, dynamical changes in atmospheric circulation are very**  
36 **uncertain<sup>4</sup>, leading to low confidence in regional projections especially for precipitation over**  
37 **the coming decades<sup>5,6</sup>. Furthermore, model simulations with tiny differences in initial condi-**  
38 **tions suggest that uncertainties may be largely irreducible due to the chaotic nature of the cli-**  
39 **mate system<sup>7-9</sup>. However, climate projections are difficult to verify until further observations**  
40 **become available. Here we assess retrospective climate predictions of the last six decades**

and show that decadal variations in north Atlantic winter atmospheric circulation are highly predictable. Crucially, climate models underestimate the predictable signal by an order of magnitude and skill is achieved despite a lack of agreement between individual model simulations. Consequently, skilful climate predictions of European and eastern North American winters are possible but require 100 times more ensemble members than would perfect models and post-processing to overcome underestimated teleconnections. Our results highlight the pressing need to understand why the signal-to-noise ratio is too small in climate models<sup>10</sup>, and the extent to which correcting this model error would reduce uncertainties in regional climate change on timescales beyond a decade.

Global climate models are used extensively to understand the drivers of past climate variability and change, and to predict what is likely to happen in the future<sup>1-3</sup>. Underpinning this is a need for accurate estimates of signals and associated uncertainties in climate model simulations in order to distinguish between different causes of past climate change, and to provide reliable confidence limits on future projections. Uncertainties are typically partitioned into three sources<sup>11</sup>: scenario uncertainty arising from an imperfect knowledge of external forcing factors, including changes in greenhouse gases, ozone, anthropogenic and volcanic aerosols, and solar irradiance; modelling uncertainty arising from the fact that different models respond differently to the same radiative forcing; and internal variability of climate that would occur in the absence of any external forcing.

Climate projections for many regions are currently highly uncertain, especially for atmospheric circulation<sup>4,12</sup> and related impacts, including precipitation<sup>5,6</sup>. This is particularly well

illustrated by the fact that modelling<sup>13,14</sup> and internal variability<sup>7,8</sup> uncertainties are each large enough to allow opposite projections of European winters, especially for the coming decades. Whilst modelling uncertainties might be reduced as models improve, internal variability uncertainties have been interpreted to be largely irreducible<sup>7-9</sup> suggesting that confident projections of European winters may never be possible. However, such conclusions assume that signals and uncertainties diagnosed from climate models are correct. Although multi-decadal and longer climate projections are difficult to verify until future observations become available, signals over the first 10 years can be more robustly evaluated using retrospective decadal predictions (hereafter referred to as hindcasts).

We use a very large multi-model ensemble of decadal hindcasts from the Coupled Model Intercomparison Project (CMIP) phases 5<sup>15</sup> and 6<sup>16</sup>. We focus on the boreal winter period (December to March) averaged over forecast years 2 to 9 to avoid seasonal to annual predictability and focus on decadal timescales. We use hindcasts starting each year over the period 1960 to 2005 from 6 CMIP5 and 8 CMIP6 modelling systems, giving a total of 169 ensemble members which are weighted equally (see Methods, Table1). Hence our total hindcast dataset comprises 77,740 (46 start dates times 169 ensemble members times 10 years) years of model integrations to provide robust statistics.

To compare with uncertainties in climate projections<sup>5,7,8,13,14</sup> we focus on European winters which are largely controlled by the North Atlantic Oscillation (NAO), the leading mode of atmospheric circulation variability in the north Atlantic<sup>17</sup>. The NAO represents the meridional gradient

in mean sea level pressure (mslp), typically measured as the difference in pressure between the Azores and Iceland. Its positive (negative) phase reflects an increased (reduced) pressure gradient driving stronger (weaker) mid-latitude westerly winds with increased (reduced) storminess, and a northward (southward) shift of the jet stream. Impacts of the NAO are characterised by a quadrupole pattern, with a positive (negative) NAO driving warmer, wetter (colder, drier) conditions in northern Europe and south-east North America along with colder, drier (warmer, wetter) conditions in southern Europe and north-east North America.

We assess skill using two different measures (see Methods): anomaly correlation coefficient (ACC) which measures the phase of variability, and mean-squared-skill-score (MSSS) which measures the amplitude of variability. We find significant skill for decadal predictions of winter mslp in most regions, including the north Atlantic, when measured by the ACC between the 169-member ensemble mean and observations (Figure 1a). However, skill is much lower especially in the north Atlantic when measured by the MSSS or the ACC of a smaller (10-member, typical of individual prediction systems<sup>16</sup>) ensemble mean (Figure 1 b and c). Timeseries from the observations and each model ensemble member consist of a predictable component (the signal) and unpredictable internal variability (the noise). The discrepancy in skill between ACC and MSSS, and the need for a large ensemble, arise because the signal-to-noise ratio is too small in the models compared to observations<sup>10, 18, 19</sup>. Hence, skill is low in a 10-member ensemble mean because a larger ensemble is required to reduce the noise and extract the predicted signal. In contrast, the signal resulting from a large ensemble mean may capture the correct phase of observed variability giving a significant ACC, but its amplitude will be much too small resulting in a low MSSS.

Errors in the signal-to-noise ratio can be quantified by comparing the predictable components (the predictable fraction of the total variability) in observations and models. The ratio of predictable components<sup>10, 18, 20</sup> (RPC, see Methods) is expected to be one for a perfect forecasting system; values greater than one show where the signal-to-noise ratio is erroneously too small in models. Consistent with differences in ACC and MSSS we find RPC is greater than one almost everywhere where there is skill in ACC, and especially in the north Atlantic (Figure 1d).

The NAO exhibits marked decadal variability<sup>21</sup> with a strong increase from the 1960s to the 1990s and a decrease thereafter (Figure 2a, black curve). The raw ensemble mean forecast shows virtually no signal (Figure 2a, red curve), and the observations generally lie within the model uncertainties (shading showing the 5-95% range diagnosed from the ensemble spread), although the extreme values in the early 1960s and late 1980s are not well-captured by models in agreement with other studies<sup>22, 23</sup>. Taken at face value, as is done for climate projections<sup>5, 7, 8, 14</sup>, the small model signal and much larger spread would imply little ability to predict the NAO and a large component of unpredictable internal variability. However, by comparing with observations we find significant correlation skill of the ensemble mean (ACC=0.48, p=0.02), while persistence provides a poor forecast (ACC=0.1). Hence, skilful climate model predictions of the NAO are possible using the ensemble mean, but the signal-to-noise ratio is too small (RPC=4.2) and its variance must be calibrated to provide realistic forecasts<sup>19</sup>.

Rescaling the ensemble mean time-series to have the same variance as the observations reveals that the predictions do capture the observed increase from the 1960s to 1990s and decrease



thereafter (Figure 2b). However, even with 169 ensemble members (Figure 2b thin red curve) there are large interannual variations that are not expected or observed in 8-year rolling means. We therefore create a larger lagged ensemble by taking the average of the four latest forecasts available at each start date (giving 676 members, Figure 2b thick red curve, see Methods). This reveals that the NAO is highly predictable on decadal timescales ( $ACC=0.79$ ,  $p<0.01$ ) in stark contrast to the lack of predictability implied by the standard interpretation of raw model output (Figure 2a). Importantly, the signal-to-noise ratio is much too small in the models ( $RPC=11$ ,  $p=0.02$ ). The total 8-year variability of the NAO in individual model members (standard deviation = 1.7 to 2.6 hPa, 5-95% range, year 2-9 hindcasts) is not significantly different to the observations (2.4hPa). Hence the predictable signal (see Methods) is underestimated by an order of magnitude in the model ensemble. Since the standard error of the ensemble mean is reduced by the square root of the ensemble size, the ensemble required to extract the signal is 100 times larger than it would be for perfect models.

The fact that the NAO signal is much too weak in models implies that the impacts of the NAO will be underestimated relative to other factors such as greenhouse gases. Hence in regions influenced by the NAO the ensemble mean will not reflect the true balance of driving factors and simply inflating its variance to be the same as observed will not correct the error. A potential solution is to post-process the model output by selecting a subset of (20) ensemble members from the lagged ensemble (of 676 members) whose simulated NAO is closest in sign and magnitude to the ensemble mean NAO after adjusting this to take into account the underestimated signal. These members contain close-to the correct magnitude of the forecast NAO whilst retaining influ-

ences from greenhouse gases and other sources. We refer to this procedure as “NAO-matching” (see Methods) and note that it builds on previous techniques<sup>24,25</sup> by using the models as much as possible instead of observed relationships which may not be causal or robust.

We investigate this technique first for forecasts of Atlantic Multidecadal Variability (AMV, see Methods). AMV is thought to be one of the most predictable aspects of decadal climate<sup>26</sup>, yet the lagged ensemble mean does not capture the correct timing of the minimum in the late 1980s (Figure 2c). NAO-matching captures the minimum and subsequent rapid warming in much better agreement with observations (Figure 2d) consistent with evidence that AMV is at least partly forced by the NAO<sup>27–29</sup>. We find similar improvements for northern European rainfall: the lagged ensemble mean is not significantly skilful and the observations lie outside the modelled uncertainties in the 1960s and 1980s (Figure 2e), whereas the NAO-matched ensemble is significantly skilful ( $ACC=0.72$ ,  $p<0.01$ ) and captures the observed increase from the 1960s to late 1980s and decrease thereafter. As expected, these improvements are not seen by simply adjusting the variance of the ensemble mean (Supplementary Figure S1).

Forecasts of extreme decades would be of particular value since they could enable action to be taken in advance to avoid the most severe climate impacts<sup>30</sup>. We therefore investigate the extreme positive NAO period between 1986 and 1997 (8-year means starting 1986 to 1990, Figure 2a). Consistent with the above results, the raw lagged ensemble mean shows virtually no signal compared to observed variability (Figure 3 a, b, c compared to d, e, f). Adjusting its variance to be equal to the observed variance (Figure 3 g, h, i) reveals that the forecasts do capture the positive

NAO (as expected from Figure 2b), but the expected impacts are underestimated, especially for temperature and northern European precipitation. However, the NAO-matched forecast (Figure 3 j, k, l) shows a clear improvement and captures the expected quadrupole pattern with warm, wet (cold, dry) anomalies in northern Eurasia and south-east North America (northern Africa and parts of southern Europe, and north-east North America), as well as low pressure across the Arctic. Similar improvements from NAO-matching are found for trends and for skill measured over all of the hindcasts (Supplementary Figures S3-S4).

We have shown that the winter NAO and related impacts on Europe and eastern North America are highly predictable on decadal timescales. AMV is usually believed to be a major source of decadal prediction skill<sup>26,31</sup>. However, we find that predictions of AMV can be improved by using the forecast NAO (Figure 2c,d), whereas predictions of the NAO are degraded by selecting the most skilful AMV ensemble members (Supplementary Figure S5). This suggests that the NAO is not solely driven by AMV. Hence other potential influences, including for example the tropics<sup>32–34</sup>, warrant further investigation.

Crucially we find that the NAO signal is underestimated by an order of magnitude in the model ensemble. This adds to an increasing body of evidence that the signal-to-noise ratio is too small in climate models, seen on seasonal<sup>20,35–37</sup>, interannual<sup>38</sup> and decadal<sup>19,39</sup> timescales. Consequently, the real world is more predictable than climate models suggest<sup>10,18</sup> and uncertainties diagnosed from raw model simulations are too large. The cause of this error is not yet known, though there are several hypotheses including weak teleconnections to the quasi-biennial

oscillation<sup>40</sup>, lack of persistence in the NAO<sup>41,42</sup> and in weather regimes<sup>43</sup>, unresolved ocean at-  
mosphere interactions<sup>44</sup> and weak transient eddy feedback<sup>45</sup>.

A key question is whether climate models also underestimate signals on timescales beyond a  
decade. There is some evidence that the atmospheric circulation response to Arctic sea ice loss<sup>46</sup>,  
and to external factors<sup>10</sup> including volcanic eruptions, solar variations and ozone changes, are too  
weak in models. Models also appear to underestimate the magnitude of multi-decadal temperature  
variability<sup>47,48</sup> especially for the north Atlantic<sup>49,50</sup>. Furthermore, model-simulated winter climate  
change signals in the north Atlantic increase substantially as resolution increases<sup>51</sup>, consistent  
with the suggestion that eddy feedbacks are inadequately resolved<sup>45</sup>. If this is robust, treating  
current model simulations at face value is giving misleading conclusions about uncertainties and  
irreducible internal variability.

## **Methods**

**Observations and models.** Near surface temperature observations are computed as the average  
of HadCRUT4<sup>52</sup>, NASA-GISS<sup>53</sup> and NCDC<sup>54</sup>. Precipitation observations are taken from GPCC<sup>55</sup>

and mean sea level pressure is taken from HadSLP2<sup>56</sup>.

We assess a large multi-model ensemble (169 members, Table 1) of decadal predictions from 14 modelling systems using hindcasts starting each year from 1960 to 2005 from the Coupled Model Intercomparison Project (CMIP5) phases 5<sup>15</sup> and 6<sup>16</sup>. We found no significant difference in NAO correlation skill between the CMIP5 and CMIP6 ensembles and focus on the combined ensemble to obtain the most robust statistics. We create ensemble means by taking the equally-weighted average of all ensemble members and assess rolling 8-year boreal winter (December to March) means defined by calendar years 2 to 9 from each start date. The forecasting systems start between 1st of November and January each year, giving a lead time of at least a year before the assessed forecast period to focus on decadal timescales and avoid predictability arising from seasonal to annual variability. Both halves of the 8-year period contribute to skill (NAO ACC = 0.57 and 0.45,  $p=0.03$ , for forecast years 2 to 5 and 6 to 9 respectively). Both observations and models were interpolated to a 5° longitude by 5° latitude grid before comparison.

**Indices.** The North Atlantic Oscillation (NAO) index is calculated as the difference in mean sea level pressure between two small boxes located around the Azores (28-20°W, 36-40°N) and Iceland (25-16°W, 63-70°N) with the average over the whole time series removed to create anomalies<sup>38</sup>. Atlantic Multidecadal Variability (AMV) is calculated as the near-surface temperature in the North Atlantic (80-0°W, 0-60°N) minus the global average (60°S-60°N)<sup>57</sup>. European rainfall is averaged over the box 10°W-25°E, 55-70°N. All forecasts indices are based on the ensemble mean.

**Forecast quality and uncertainty measures.** Model biases and drifts are treated by computing anomalies relative to climatology for each model computed over all hindcasts, and comparing with observed anomalies computed over the same period. Although there are many ways to measure forecast quality, we focus on those that illustrate the underestimated model signals by using the following:

$$\text{Pearson anomaly correlation coefficient ACC} = \frac{\sum_{i=1}^N (f_i - \bar{f})(o_i - \bar{o})}{\sum_{i=1}^N (f_i - \bar{f})^2 \sum_{i=1}^N (o_i - \bar{o})^2} \quad (1)$$

$$\text{Mean-squared-skill-score MSSS} = 1 - \frac{\sum_{i=1}^N (f_i - o_i)^2}{\sum_{i=1}^N (\bar{o} - o_i)^2} \quad (2)$$

$$\text{Ratio of predictable components RPC} = \frac{\sigma_{sig}^o / \sigma_{tot}^o}{\sigma_{sig}^f / \sigma_{tot}^f} = \frac{ACC}{\sigma_{sig}^f / \sigma_{tot}^f} \quad (3)$$

$$\text{Ratio of predictable signals} = \frac{\sigma_{sig}^o}{\sigma_{sig}^f} = RPC \frac{\sigma_{tot}^o}{\sigma_{tot}^f} \quad (4)$$

where  $N$  is the number of hindcast start dates,  $f_i$  and  $o_i$  are the ensemble mean forecast and observations at each time, and the overbar represents the average over all times.  $\sigma_{sig}$  and  $\sigma_{tot}$  are the expected standard deviations of the predictable signal and total variability, with superscripts o and f for the observations and forecasts respectively. For the forecasts,  $\sigma_{sig}$  and  $\sigma_{tot}$  are computed from the ensemble mean and individual members respectively.

ACC measures the ability to predict the phase of variability, whereas MSSS measures the magnitude of errors relative to a climatological forecast. For a perfect forecasting system RPC should equal one. Note that RPC is not computed where the ACC is negative, and that the above formula likely gives a lower bound<sup>10,18</sup>.

Uncertainties in raw model forecasts are computed from the ensemble standard deviation

for each start date. Uncertainties in variance adjusted and NAO-matched forecasts are computed from the root-mean-square error between the ensemble mean and the observations as required for reliable forecasts<sup>58</sup>.

We note that it is theoretically possible for the multi-model RPC to be larger than for individual models if time dependent model biases<sup>59</sup> or teleconnection errors reduce the model signal more than the correlation with observations. Assessing this thoroughly would require large ensembles of individual model hindcasts which are not available. However, assessing the largest individual model ensemble available (NCAR CESM1.1 with 40 members per year, giving 160 lagged members, Table 1) does not support this hypothesis: the NCAR RPC of 6.2 is not significantly different from the average RPC of multi-model ensembles of the same size (4.8 averaged over 1000 random samples, with 5-95% range 1.3 to 7.4). Furthermore, the statistics presented in this study are appropriate for multi-model ensemble forecasts.

We further note that there is some evidence that the predictability of the NAO may vary on multi-decadal timescales<sup>60</sup>, though this is not robust across models<sup>61</sup>. Our results are statistically significant for the hindcast period available, but longer hindcasts that include more cycles of decadal variability would be beneficial for future studies.

**Lagged ensemble.** Consecutive 8-year means contain 7 identical years. Hence large interannual variations, as seen in 169-member ensemble mean NAO forecasts (Figure 2b), are not expected. They occur because the signal to noise ratio is too small in models and consecutive decadal predictions consist of independent model simulations that are dominated by different samples of the

noise. Ideally additional ensemble members would be used to reduce the noise further, but these are not available. Instead we create a lagged ensemble by combining the required forecast with the previous three i.e. the year 2-9 forecasts starting in 1963 are combined with the year 2-9 forecasts starting in 1962, 1961 and 1960 giving a total of 676 members (169 members time 4 start dates). The previous forecasts are sub-optimal because they do not cover exactly the same forecast period, and rely on the persistence of running 8-year means. Hence there is a trade off between reducing the noise with additional members and potentially degrading the skill by relying on persistence. In the current generation of climate models the benefit in reducing the noise far outweighs the degradation from using persistence. We present results for the combination of 4 lagged forecasts, but find similar levels of skill for other combinations (NAO ACC = 0.71 and 0.78 for combining 3 and 5 lagged forecasts respectively). A similar technique relying on persistence of the predictor recently proved to strongly reduce the noise in decadal predictions of summer temperature extremes over land<sup>62</sup>.

**NAO-matching.** At any location that is influenced by the NAO we can write

$$O = O_{NAO} + O_{OTHER} + \epsilon^o \quad (5)$$

$$F^k = F_{NAO}^k + F_{OTHER}^k + \epsilon^k \quad (6)$$

$$\hat{F} = \hat{F}_{NAO} + \hat{F}_{OTHER} + \hat{\epsilon} \quad (7)$$

where  $O$ ,  $F^k$  and  $\hat{F}$  are the observed, forecast ensemble member  $k$  and forecast ensemble mean values of a meteorological variable (e.g. temperature, rainfall, pressure). The subscript *NAO* refers to the portion that is related to the NAO, the subscript *OTHER* refers to the portion related to other predictable drivers (including greenhouse gases and sea surface temperatures unrelated to



the NAO) and  $\epsilon$  is an unpredictable residual. Because the predictable NAO signal is too small in models, the mean of a very large ensemble is required for skilful NAO predictions (Figure 2b). However, the magnitude of the ensemble mean NAO is much too small (Figure 2a) and therefore  $\hat{F}_{NAO}$  will be severely underestimated.

One approach to overcoming model deficiencies uses regressions between model hindcasts and observations<sup>25, 63–65</sup>, which effectively replaces the erroneous  $\hat{F}_{NAO}$  with the observed value  $O_{NAO}$ . Whilst this can give very good results, it relies on  $O_{NAO}$  estimated from the observations being robust and describing a causal relationship between the NAO and remote regions. This approach is less attractive on decadal than seasonal timescales because  $O_{NAO}$  is potentially more affected by sampling errors from the relatively small hindcast period.

An alternative approach<sup>24</sup> replaces the underestimated  $\hat{F}_{NAO}$  with more realistic  $F_{NAO}^k$  by selecting from the full ensemble a smaller set of members that have the required magnitude of the NAO. These members contain close-to the correct magnitude of the required NAO and its teleconnections whilst retaining other influences. Hence,  $\hat{F}_{NAO}$  for this selected ensemble will be larger than that of the full ensemble, thereby increasing the signal. Because the selected ensemble is smaller the remaining noise will not be reduced as much as in the full ensemble. However, the selection process transfers variability from what would be considered as noise in a random ensemble into  $\hat{F}_{NAO}$ , thereby reducing  $\hat{\epsilon}$  in the selected ensemble. Hence, in regions affected by the NAO the increase in signal is likely to be larger than the reduced suppression of the remaining noise, thereby increasing the signal to noise ratio and improving the skill.

In the previous seasonal forecast study<sup>24</sup> the required NAO was obtained based on observed relationships with potential drivers. However, on decadal timescales such relationships are not well-established and are more likely to be affected by sampling errors. We therefore take the required NAO to be the ensemble mean forecast NAO but adjusted to account for the underestimation of the predictable signal. This is achieved by multiplying the ensemble mean NAO by the ratio of predictable signals (equation 4). To avoid overfitting to observations we compute the ratio of predictable signals for each hindcast start date separately using a cross-validation approach in which the required hindcast and those on either side are omitted. Our conclusions are robust to omitting more hindcasts (we have tested up to 4 years either side) though skill may be underestimated especially in these cases<sup>66,67</sup>.

The overall procedure is as follows. For each start date  $i$ :

1. Compute the signal-adjusted (described above) NAO index of the ensemble mean  $\hat{\text{NAO}}_i$
2. Compute the NAO index for each ensemble member  $\text{NAO}_i^k$
3. For each ensemble member calculate the difference  $\text{NAO}_i^k - \hat{\text{NAO}}_i$
4. Select the  $M$  ( $= 20$ ) ensemble members with the smallest absolute differences

We take the mean of this subset of  $M$  members and present standardised forecast anomalies (Figure 3) or adjust its variance to be the same as observed (Figure 2). We note that this approach is applicable to forecasts as well as hindcasts. We present results for a subset of 20 members, but the results are similar for subsets ranging from 10 to 40 members. This method relies on models

simulating realistic NAO teleconnections ( $F_{NAO}^k$ ) and further improvements might be possible by using the best models in this respect, but this is beyond the scope of this study.

**Significance.** For a given set of validation cases, we test for values that are unlikely to be accounted for by uncertainties arising from a finite ensemble size ( $E$ ) and a finite number of validation points ( $N$ ). This is achieved using a non-parametric block bootstrap approach<sup>19,68,69</sup>, in which an additional 1000 hindcasts are created as follows:

1. Randomly sample with replacement  $N$  validation cases. In order to take autocorrelation into account this is done in blocks of 5 consecutive cases.
2. For each of these, randomly sample with replacement  $E$  ensemble members.
3. Compute the required statistic for the ensemble mean (e.g. correlation, MSSS, RPC).
4. Repeat from (1) 1000 times to create a probability distribution.
5. Obtain the significance level based on a 2-tailed test of the hypothesis that skill is zero, or RPC is one.

**Data Availability** The datasets analysed during the current study are available from the CMIP data archives.

**Code Availability** The code used during the current study is available from the corresponding author on reasonable request.

**Acknowledgements** DMS, AAS, NJD, LH and RE were supported by the Met Office Hadley Centre Climate Programme funded by BEIS and Defra and by the European Commission Horizon 2020 EUCP

project (GA 776613). FJDR, LPC, SW and RB also acknowledge the support from the EUCP project (GA 776613) and from the Ministerio de Economía y Competitividad (MINECO) as part of the CLINSA project (Grant No. CGL2017-85791-R). SW received funding from the innovation programme under the Marie Skłodowska-Curie grant agreement H2020-MSCA-COFUND-2016-754433 and PO from the Ramon y Cajal senior tenure programme of MINECO. The EC-Earth simulations were performed on Marenostrum 4 (hosted by the Barcelona Supercomputing Center, Spain) using Auto-Submit through computing hours provided by PRACE. WAM, HP, KM and KP were supported by the German Federal Ministry for Education and Research (BMBF) project MiKlip (grant 01LP1519A). NK, IB, FC and YW were supported by the Norwegian Research Council projects SFE (grant 270733) the Nordic Center of excellent ARCPATH (grant 76654) and the Trond Mohn Foundation, under the project number : BFS2018TMT01 and received grants for computer time from the Norwegian Program for supercomputing (NOTUR2, NN9039K) and storage grants (NORSTORE, NS9039K). JM, LFB and DS are supported by Blue-Action (European Union Horizon 2020 research and innovation program, Grant Number: 727852) and EUCP (European Union Horizon 2020 research and innovation programme under grant agreement no 776613) projects. The National Center for Atmospheric Research (NCAR) is a major facility sponsored by the US National Science Foundation (NSF) under Cooperative Agreement No. 1852977. NCAR contribution was partially supported by the National Oceanic and Atmospheric Administration (NOAA) Climate Program Office under Climate Variability and Predictability Program Grant NA13OAR4310138 and by the US NSF Collaborative Research EaSM2 Grant OCE-1243015.

**Competing Interests** The authors declare that there are no competing interests.

**Correspondence** Correspondence and requests for materials should be addressed to D.M.S. (email: doug.smith@metoffice.gov.uk).

**Author contributions** D.M.S. led the analysis and writing with comments from all authors. R.E. processed the CMIP5 data. A.A.S. suggested NAO-matching. All authors except A.A.S., P.A., A.B., P.-A.M., D.N., J.R. and P.R. contributed to creating the decadal prediction data.

1. Bindoff, N. L. *et al.* Detection and attribution of climate change: from global to regional. In Stocker, T. F. *et al.* (eds.) *Climate Change 2013: The Physical Science Basis. Contribution of Working Group I to the Fifth Assessment Report of the Intergovernmental Panel on Climate Change* (Cambridge University Press, 2013).
2. Kirtman, B. *et al.* Near-term climate change: Projections and predictability. In Stocker, T. F. *et al.* (eds.) *Climate Change 2013: The Physical Science Basis. Contribution of Working Group I. to the Fifth Assessment Report of the Intergovernmental Panel on Climate Change* (Cambridge University Press, 2013).
3. Collins, M. *et al.* Long-term climate change: Projections, commitments and irreversibility. In Stocker, T. F. *et al.* (eds.) *Climate Change 2013: The Physical Science Basis. Contribution of Working Group I to the Fifth Assessment Report of the Intergovernmental Panel on Climate Change*, 1029–1136 (Cambridge University Press, 2013).
4. Shepherd, T. G. Atmospheric circulation as a source of uncertainty in climate change projections. *Nature Geosci.* **7**, 703–708 (2014).
5. Hawkins, E. & Sutton, R. The potential to narrow uncertainty in projections of regional precipitation change. *Clim. Dyn.* **37**, 407–418 (2011).

6. Knutti, R. & Sedleck, J. Robustness and uncertainties in the new CMIP5 climate model projections. *Nature Climate Change* **3**, 369–373 (2013).
7. Hawkins, E., Smith, R. S., Gregory, J. M. & Stainforth, D. A. Irreducible uncertainty in near-term climate projections. *Clim. Dyn.* **46**, 3807–3819 (2016).
8. Deser, C., Hurrell, J. W. & Phillips, A. S. The role of the North Atlantic Oscillation in European climate projections. *Clim. Dyn.* **49**, 3141–3157 (2017).
9. Marotzke, J. Quantifying the irreducible uncertainty in near term climate projections. *Wiley Interdisciplinary Reviews: Climate Change* **10**, e563 (2019).
10. Scaife, A. A. & Smith, D. A signal-to-noise paradox in climate science. *npj Climate and Atmospheric Science* **1**, 28 (2018).
11. Hawkins, E. & Sutton, R. The Potential to Narrow Uncertainty in Regional Climate Predictions. *Bull. Am. Meteorol. Soc.* **90**, 1095–1108 (2009).
12. Fereday, D., Chadwick, R., Knight, J. & Scaife, A. A. Atmospheric Dynamics is the Largest Source of Uncertainty in Future Winter European Rainfall. *J. Climate* **31**, 963–977 (2018).
13. Woollings, T. Dynamical influences on european climate: an uncertain future. *Philos. Trans. R. Soc. London* **368**, 3733–3756 (2010).
14. Zappa, G. & Shepherd, T. G. Storylines of Atmospheric Circulation Change for European Regional Climate Impact Assessment. *J. Climate* **30**, 6561–6577 (2017).

- 389 15. Taylor, K. E., Stouffer, R. J. & Meehl, G. A. An overview of CMIP5 and the experiment  
390 design. *Bull. Am. Meteorol. Soc.* **93**, 485–498 (2012).
- 391 16. Boer, G. J. *et al.* The Decadal Climate Prediction Project (DCPP) contribution to CMIP6.  
392 *Geosci. Model Devel.* (2016).
- 393 17. Hurrell, J. W., Kushnir, Y., Ottersen, G. & Visbeck, M. (eds.) *The North Atlantic Oscillation:  
394 Climatic Significance and Environmental Impact*, vol. 134 of *Geophysical Monograph Series*  
395 (American Geophysical Union, Washington, D. C., 2003).
- 396 18. Eade, R. *et al.* Do seasonal-to-decadal climate predictions underestimate the predictability of  
397 the real world? *Geophys. Res. Lett.* **41**, 5620–5628 (2014).
- 398 19. Smith, D. M. *et al.* Robust skill of decadal climate predictions. *npj Climate and Atmospheric  
399 Science* **2**, 13 (2019).
- 400 20. Siegert, S. *et al.* A Bayesian framework for verification and recalibration of ensemble fore-  
401 casts: How uncertain is NAO predictability? *J. Climate* **29**, 995–1012 (2016).
- 402 21. Hurrell, J. W. Decadal trends in the North Atlantic Oscillation: regional temperatures and  
403 precipitation. *Science* **269**, 676–679 (1995).
- 404 22. Scaife, A. A. *et al.* The CLIVAR C20C project: selected twentieth century climate events.  
405 *Clim. Dyn.* **33**, 603–614 (2009).
- 406 23. Bracegirdle, T. J., Lu, H., Eade, R. & Woollings, T. Do CMIP5 Models Reproduce Observed  
407 Low Frequency North Atlantic Jet Variability? *Geophys. Res. Lett.* **45**, 7204–7212 (2018).

24. Dobrynin, M. *et al.* Improved Teleconnection-Based Dynamical Seasonal Predictions of Boreal Winter. *Geophys. Res. Lett.* **45**, 3605–3614 (2018).
25. Simpson, I. R., Yeager, S. G., McKinnon, K. A. & Deser, C. Decadal predictability of late winter precipitation in western Europe through an oceanjet stream connection. *Nature Geosci.* **12**, 613–619 (2019).
26. Yeager, S. G. & Robson, J. I. Recent progress in understanding and predicting Atlantic decadal climate variability. *Current Climate Change Reports* **3**, 112–127 (2017).
27. Eden, C. & Willebrand, J. Mechanism of interannual to decadal variability of the North Atlantic circulation. *J. Climate* **14**, 2266–2280 (2001).
28. McCarthy, G. D., Haigh, I. D., Hirschi, J. J.-M., Grist, J. P. & Smeed, D. A. Ocean impact on decadal Atlantic climate variability revealed by sea-level observations. *Nature* **521**, 508–510 (2015).
29. Clement, A. *et al.* The Atlantic Multidecadal Oscillation without a role for ocean circulation. *Science* **350**, 320–324 (2015).
30. Zanardo, S., Nicotina, L., Hilberts, A. G. J. & Jewson, S. P. Modulation of Economic Losses From European Floods by the North Atlantic Oscillation. *Geophys. Res. Lett.* **46**, 2563–2572 (2019).
31. Eden, C., Greatbatch, R. J. & Lu, J. Prospects for decadal prediction of the North Atlantic Oscillation (NAO). *Geophys. Res. Lett.* **29**, 104–1–104–4 (2002).



- 427 32. Hoerling, M. P., Hurrell, J. W. & Xu, T. Tropical origins for recent North Atlantic climate  
428 change. *Science* **292**, 90–92 (2001).
- 429 33. Greatbatch, R. J., Lin, H., Lu, J., Peterson, K. A. & Derome, J. Tropical/Extratropical forcing  
430 of the AO/NAO: A corrigendum. *Geophys. Res. Lett.* **30** (2003).
- 431 34. Shin, S.-I. & Sardeshmukh, P. D. Critical influence of the pattern of Tropical Ocean warming  
432 on remote climate trends. *Clim. Dyn.* **36**, 1577–1591 (2011).
- 433 35. Scaife, A. A. *et al.* Skillful long-range prediction of european and north american winters.  
434 *Geophys. Res. Lett.* **41**, 2514–2519 (2014).
- 435 36. Dunstone, N. J. *et al.* Skilful seasonal predictions of summer European rainfall. *Geophys.*  
436 *Res. Lett.* (2018).
- 437 37. Baker, L. H., Shaffrey, L. C., Sutton, R. T., Weisheimer, A. & Scaife, A. A. An intercomparison  
438 of skill and over/underconfidence of the wintertime North Atlantic Oscillation in multi-model  
439 seasonal forecasts. *Geophys. Res. Lett.* (2018).
- 440 38. Dunstone, N. J. *et al.* Skilful predictions of the winter North Atlantic Oscillation one year  
441 ahead. *Nature Geosci.* (2016).
- 442 39. Yeager, S. G. *et al.* Predicting near-term changes in the earth system: A large ensemble of  
443 initialized decadal prediction simulations using the Community Earth System Model. *Bull.*  
444 *Am. Meteorol. Soc.* **99**, 1867–1886 (2018).

- 445 40. O'Reilly, C. H., Weisheimer, A., Woollings, T., Gray, L. J. & MacLeod, D. The importance of  
446 stratospheric initial conditions for winter North Atlantic Oscillation predictability and impli-  
447 cations for the signal-to-noise paradox. *Q. J. R. Meteorol. Soc.* **145**, 131–146 (2019).
- 448 41. Zhang, W. & Kirtman, B. Understanding the Signal-to-Noise Paradox with a Simple Markov  
449 Model. *Geophys. Res. Lett.* 2019GL085159 (2019).
- 450 42. Jin, Y., Rong, X. & Liu, Z. Potential predictability and forecast skill in ensemble climate  
451 forecast: a skill-persistence rule. *Clim. Dyn.* **51**, 2725–2741 (2018).
- 452 43. Strommen, K. & Palmer, T. N. Signal and noise in regime systems: A hypothesis on the  
453 predictability of the North Atlantic Oscillation. *Q. J. R. Meteorol. Soc.* **145**, 147–163 (2019).
- 454 44. Czaja, A., Frankignoul, C., Minobe, S. & Vannière, B. Simulating the Midlatitude Atmo-  
455 spheric Circulation: What Might We Gain From High-Resolution Modeling of Air-Sea Inter-  
456 actions? *Current Climate Change Reports* **5**, 390–406 (2019).
- 457 45. Scaife, A. A. *et al.* Does increased atmospheric resolution improve seasonal climate predic-  
458 tions? *Atmos. Sci. Lett.* **20** (2019).
- 459 46. Mori, M., Kosaka, Y., Watanabe, M., Nakamura, H. & Kimoto, M. A reconciled estimate  
460 of the influence of Arctic sea-ice loss on recent Eurasian cooling. *Nature Climate Change* **9**,  
461 123–129 (2019).
- 462 47. Cheung, A. H. *et al.* Comparison of Low-Frequency Internal Climate Variability in CMIP5  
463 Models and Observations. *J. Climate* **30**, 4763–4776 (2017).

48. Kravtsov, S. Pronounced differences between observed and CMIP5-simulated multidecadal climate variability in the twentieth century. *Geophys. Res. Lett.* **44**, 5749–5757 (2017).
49. Wang, X., Li, J., Sun, C. & Liu, T. NAO and its relationship with the Northern Hemisphere mean surface temperature in CMIP5 simulations. *J. Geophys. Res.* **122**, 4202–4227 (2017).
50. Kim, W. M., Yeager, S. G. & Danabasoglu, G. Key role of internal ocean dynamics in Atlantic multidecadal variability during the last half century. *Geophys. Res. Lett.* **45** (2018).
51. Baker, A. J. *et al.* Enhanced Climate Change Response of Wintertime North Atlantic Circulation, Cyclonic Activity, and Precipitation in a 25-km-Resolution Global Atmospheric Model. *J. Climate* **32**, 7763–7781 (2019).
52. Morice, C. P., Kennedy, J. J., Rayner, N. A. & Jones, P. D. Quantifying uncertainties in global and regional temperature change using an ensemble of observational estimates: The HadCRUT4 data set. *J. Geophys. Res.* **117**, D08101 (2012).
53. Hansen, J., Ruedy, R., Sato, M. & Lo, K. Global surface temperature change. *Rev. Geophys.* **48** (2010).
54. Karl, T. R. *et al.* Possible artifacts of data biases in the recent global surface warming hiatus. *Science* **348**, 1469–1472 (2015).
55. Schneider, U. *et al.* GPCC’s new land surface precipitation climatology based on quality-controlled in situ data and its role in quantifying the global water cycle. *Theor. Appl. Climatol.* **115**, 15–40 (2014).

56. Allan, R. J. & Ansell, T. J. A new globally complete monthly historical gridded mean sea level pressure data set (HadSLP2): 1850-2003. *J. Climate* **19**, 5816–5842 (2006).
57. Trenberth, K. E. & Shea, D. J. Atlantic hurricanes and natural variability in 2005. *Geophys. Res. Lett.* **33**, L12704 (2006).
58. Doblas-Reyes, F. J. *et al.* Addressing model uncertainty in seasonal and annual dynamical ensemble forecasts. *Q. J. R. Meteorol. Soc.* **135**, 1538–1559 (2009).
59. Hodson, D. L. R. & Sutton, R. T. Exploring multi-model atmospheric GCM ensembles with ANOVA. *Climate Dynamics* **31**, 973–986 (2008).
60. Weisheimer, A. *et al.* How confident are predictability estimates of the winter North Atlantic Oscillation? *Q. J. R. Meteorol. Soc.* **145**, 140–159 (2019).
61. Kumar, A. & Chen, M. Causes of skill in seasonal predictions of the Arctic Oscillation. *Climate Dynamics* **51**, 2397–2411 (2018).
62. Borchert, L. F. *et al.* Decadal predictions of the probability of occurrence for warm summer temperature extremes. *Geophys. Res. Lett.* (2019).
63. Krishnamurti, T. N. *et al.* Improved weather and seasonal climate forecasts from multimodel superensemble. *Science* **285**, 1548–1550 (1999).
64. Yun, W. T., Stefanova, L. & Krishnamurti, T. N. Improvement of the multimodel superensemble technique for seasonal forecasts. *J. Climate* **16**, 3834–3840 (2003).

65. Kug, J.-S., Lee, J.-Y., Kang, I.-S., Wang, B. & Park, C.-K. Optimal multi-model ensemble method in seasonal prediction. *Asia-Pacific Journal of Atmospheric Sciences* **44**, 259–267 (2008).
66. Gangsto, R., Weigel, A. P., Lineger, M. A. & Appenzeller, C. Methodological aspects of the validation of decadal predictions. *Climate Res.* **55**, 181–200 (2013).
67. Smith, D., Eade, R. & Pohlmann, H. A comparison of full-field and anomaly initialization for seasonal to decadal climate prediction. *Clim. Dyn.* **41**, 3325–3338 (2013).
68. Wilks, D. S. *Statistical methods in the atmospheric sciences*, vol. 100 of *International geophysics series* (Academic Press, 2011), third edn.
69. Goddard, L. *et al.* A verification framework for interannual-to-decadal predictions experiments. *Clim. Dyn.* **40**, 245–272 (2013).
70. Doblas-Reyes, F. J. *et al.* Using EC-Earth for climate prediction research. In *ECMWF Newsletter* (ECMWF, 2018).
71. Haarsma, R. *et al.* HighResMIP versions of EC-Earth: EC-Earth3P and EC-Earth3P-HR. Description, model performance, data handling and validation. *Geosci. Model Dev.* (submitted).
72. Counillon, F. *et al.* Flow-dependent assimilation of sea surface temperature in isopycnal coordinates with the Norwegian Climate Prediction Model. *Tellus A* **68**, 32437 (2016).
73. Wang, Y. *et al.* Optimising assimilation of hydrographic profiles into isopycnal ocean models with ensemble data assimilation. *Ocean Modelling* **114**, 33–44 (2017).

74. Kharin, V. V., Boer, G. J., Merryfield, W. J., Scinocca, J. F. & Lee, W.-S. Statistical adjustment of decadal predictions in a changing climate. *Geophys. Res. Lett.* **39**, L19705 (2012).
75. Swart, N. C. *et al.* The Canadian Earth System Model version 5 (CanESM5.0.3). *Geosci. Model Devel.* **12**, 4823–4873 (2019).
76. Sospedra-Alfonso, R. & Boer, G. J. Assessing the impact of initialization on decadal prediction skill. *Geophys. Res. Lett.* (2020).
77. Yang, X. *et al.* A predictable amo-like pattern in GFDL's fully-coupled ensemble initialization and decadal forecasting system. *J. Climate* **26**, 650–661 (2013).
78. Williams, K. D. *et al.* The Met Office Global Coupled model 3.0 and 3.1 (GC3.0 and GC3.1) configurations. *J. Adv. Model Earth Syst.* **10**, 357–380 (2018).
79. Müller, W. A. *et al.* Forecast skill of multi-year seasonal means in the decadal prediction system of the Max Planck Institute for Meteorology. *Geophys. Res. Lett.* **39**, L22707 (2012).
80. Pohlmann, H. *et al.* Realistic Quasi-Biennial Oscillation Variability in Historical and Decadal Hindcast Simulations Using CMIP6 Forcing. *Geophys. Res. Lett.* 2019GL084878 (2019).
81. Chikamoto, Y. *et al.* An overview of decadal climate predictability in a multi-model ensemble by climate model MIROC. *Clim. Dyn.* **40**, 1201–1222 (2012).
82. Mochizuki, T. *et al.* Decadal prediction using a recent series of MIROC global climate models. *J. Meteorol. Soc. Jpn* **90**, 373–383 (2012).

Table 1: Forecast systems and ensemble sizes.

Forecast Centre	Model	Atmosphere resolution <sup>1</sup>	Ocean resolution <sup>2</sup>	Ensemble size	CMIP version
Barcelona Supercomputing Center, Spain	EC-Earth3 <sup>70, 71</sup>	0.7x0.7x91x0.01	1x1x0.3x75	10	CMIP6
Bjerknes Center for Climate Research, Norway	NorCPM1 <sup>72, 73</sup>	1.9x2.5x26x3	0.7x1.125x0.25x53	20	CMIP6
Canadian Centre for Climate Modelling and Analysis, Environment and Climate Change Canada	CanCM4 <sup>74</sup>	2.8x2.8x35x1	0.94x1.41x40	10	CMIP5
	CanESM5 <sup>75, 76</sup>	2.8x2.8x49x1	1x1x0.3x45	10	CMIP6
Geophysical Fluid Dynamics Laboratory, USA	CM2.1 <sup>77</sup>	2x2.5x24x3	1x1x0.3x50	10	CMIP5
IPSL-EPOC, France	IPSL-CM6A-LR	1.25x2.5x79x0.005	1x1x0.3x75	10	CMIP6
Met Office Hadley Centre, UK	HadCM3 <sup>67</sup>	2.5x3.75x19x4.5	1.25x1.25x20	20	CMIP5
	HadGEM3 <sup>78</sup>	0.55x0.83x85x0.005	0.25x0.25x75	10	CMIP6
Max Planck Institute for Meteorology, Germany	MPI-ESM1.0-LR <sup>79</sup>	1.9x1.9x47x0.01	1.5x1.5x40	3	CMIP5
	MPI-ESM1.2-HR <sup>80</sup>	0.9x0.9x95x0.01	0.4x0.4x40	10	CMIP6
National Center for Atmospheric Research, USA	CESM1.1 <sup>39</sup>	0.9x1.25x30x2.26	1x1.125x0.27x60	40	CMIP6
University of Tokyo, National Institute for Environmental Studies, and Japan Agency for Marine-Earth Science and Technology, Japan	MIROC5 <sup>81, 82</sup>	1.4x1.4x40x3	1.4x1.4x0.5x49	6	CMIP5
	MIROC6	1.4x1.4x81x0.004	1x1x0.5x62	10	CMIP6

<sup>1</sup> Atmosphere resolution (degrees latitude)x(degrees longitude)x(number of vertical levels)x(lid height, hPa)

<sup>2</sup> Ocean resolution (degrees latitude)x(degrees longitude)x(optional degrees latitude at Equator)x(number of vertical levels)

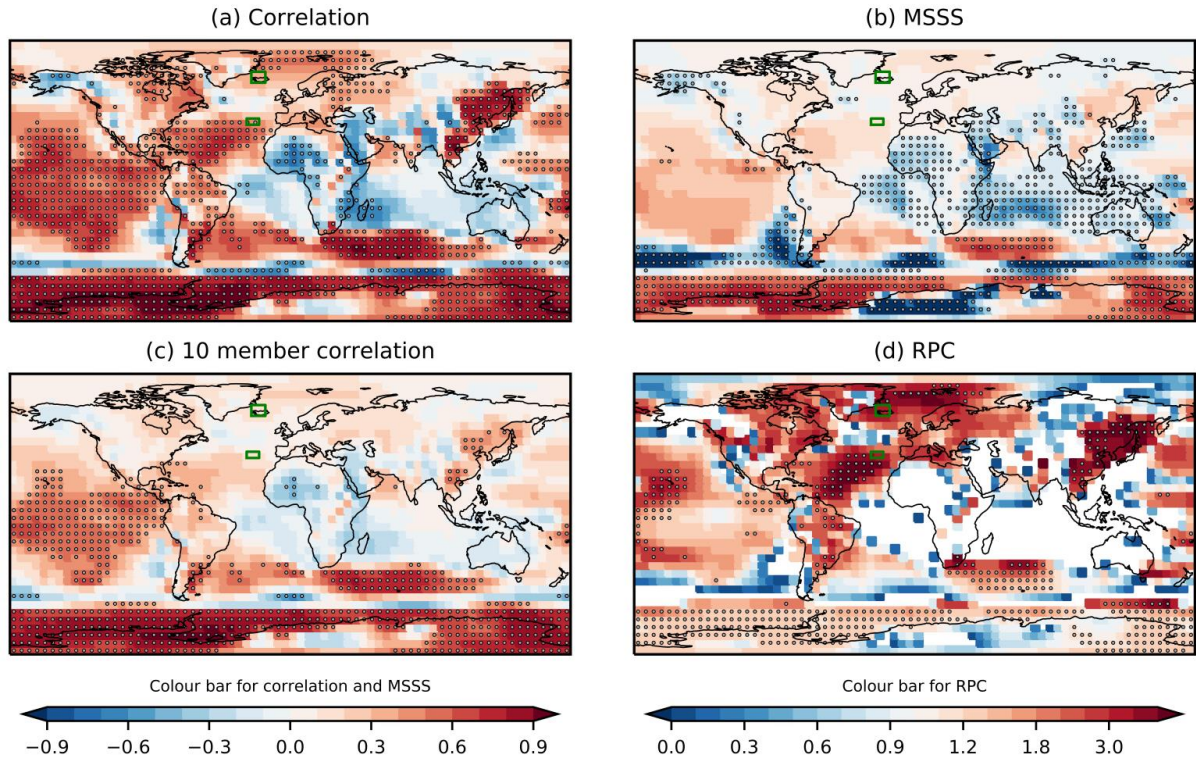
**Figure 1** Decadal prediction skill for boreal winter (December to March) mean sea level pressure. Skill for year 2-9 multi-model ensemble mean forecasts measured by (a) anomaly correlation, (b) mean squared skill score (MSSS), (c) average anomaly correlation for a 10-member ensemble mean (computed over 1000 random samples). (d) The ratio of predictable components (RPC). RPC is not calculated where the correlation is negative. Stippling shows where correlations and MSSS, or RPC, are significantly different to zero, or greater than one, respectively (95% confidence interval, see Methods). Green boxes show the regions used to calculate the NAO.

**Figure 2** Underestimated signals. (a) Time series of observed (black curve) and model forecast (years 2-9, red curve showing ensemble mean of 169 members and red shading showing the 5-95% confidence interval diagnosed from the individual members) 8-year running mean December to March NAO index. (b) As (a) but for ensemble mean forecast rescaled to have the same variance as the observations (thin red curve), and additionally smoothed by taking the lagged average of the latest four forecasts at each start date (thick red curve, 676 members, see Methods). Forecast uncertainty (red shading, 5-95% confidence interval) is obtained from the forecast ensemble mean error variance (see Methods). (c) As (a) but for AMV and lagged ensemble. (d) As (c) but for NAO-matched forecast (see Methods). (e, f) As (c, d) but for northern European rainfall. Values of anomaly correlation (ACC) of the forecast ensemble mean and of persisting the latest

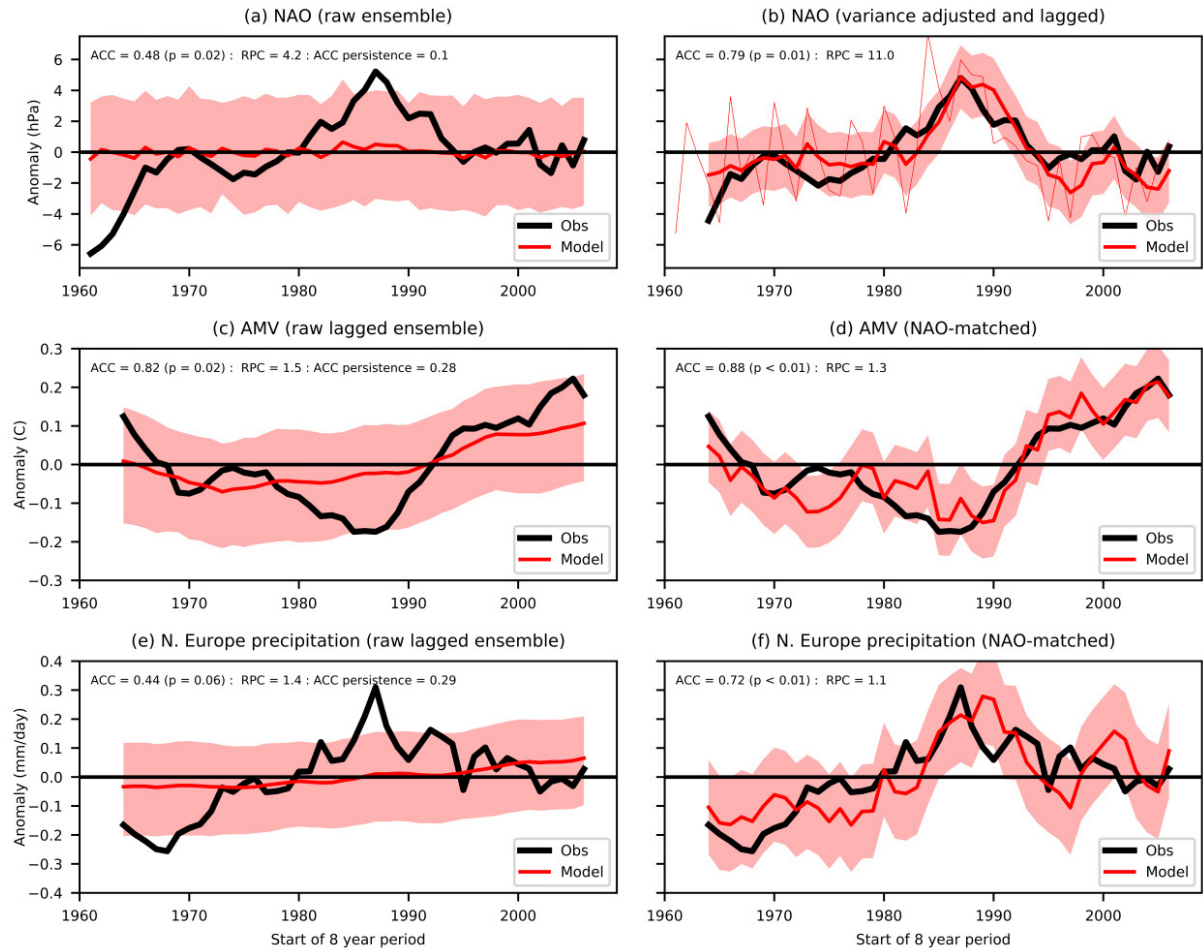


557 observed 8-year mean available before each start date, and the ratio of predictable com-  
558 ponents (RPC), are indicated. Indices are defined in Methods. Time-series are anomalies  
559 relative to the average of all year 2-9 hindcasts.

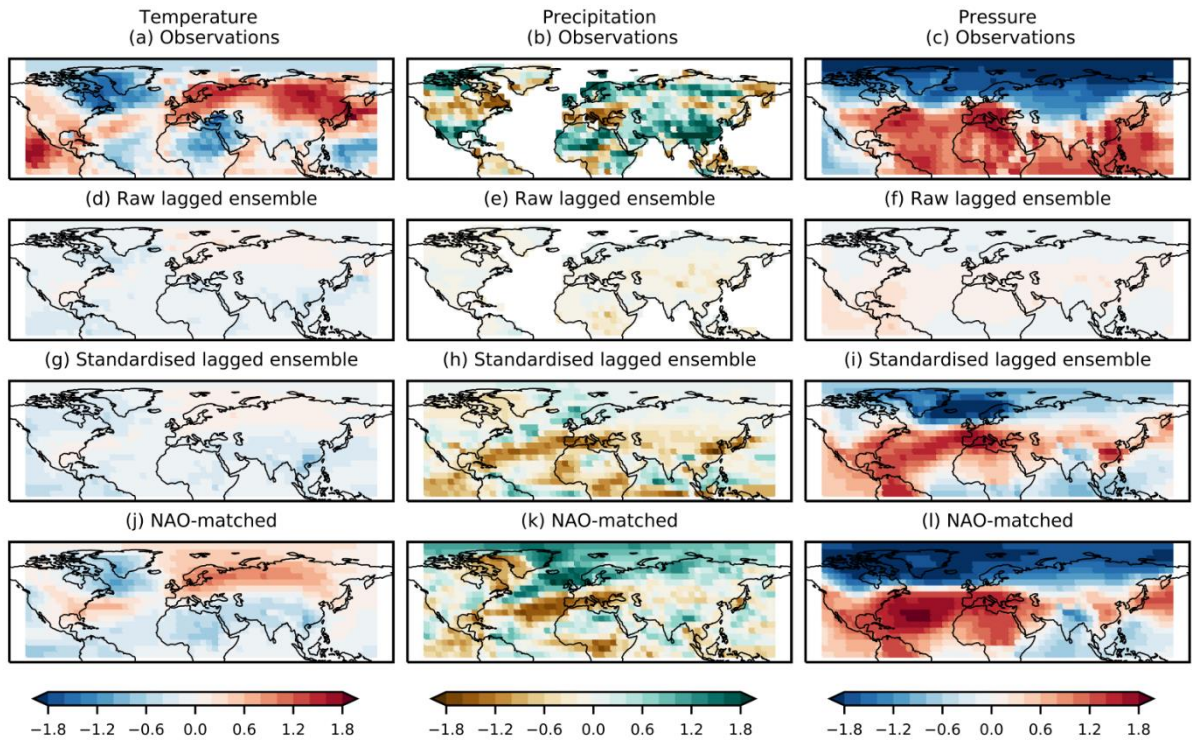
560 **Figure 3** Decadal predictions of the extreme NAO period (1986 to 1997). Observed  
561 anomalies of (a) temperature, (b) precipitation and (c) mean sea level pressure. (d, e,  
562 f) As (a, b, c) but for raw lagged ensemble mean forecasts. (g, h, i) As (d, e, f) but  
563 standardised by the ensemble mean standard deviation. (j, k, l) As (d, e, f) but for NAO-  
564 matched forecasts. Averages are taken for boreal winter (December to March) for all year  
565 2-9 forecasts verifying in the period 1986 to 1997 (i.e. start dates 1985 to 1989 inclusive),  
566 and converted to anomalies by removing the average over all hindcasts (i.e. start dates  
567 1960 to 2005 inclusive). Units are standard deviations. The raw lagged ensemble (d, e,  
568 f) is divided by the observed standard deviation to show the signal relative to observed  
569 variability.



**Figure 1: Decadal prediction skill for boreal winter (December to March) mean sea level pressure.** Skill for year 2-9 multi-model ensemble mean forecasts measured by (a) anomaly correlation, (b) mean squared skill score (MSSS), (c) average anomaly correlation for a 10-member ensemble mean (computed over 1000 random samples). (d) The ratio of predictable components (RPC). RPC is not calculated where the correlation is negative. Stippling shows where correlations and MSSS, or RPC, are significantly different to zero, or greater than one, respectively (95% confidence interval, see Methods). Green boxes show the regions used to calculate the NAO.



**Figure 2: Underestimated signals.** (a) Time series of observed (black curve) and model forecast (years 2-9, red curve showing ensemble mean of 169 members and red shading showing the 5-95% confidence interval diagnosed from the individual members) 8-year running mean December to March NAO index. (b) As (a) but for ensemble mean forecast rescaled to have the same variance as the observations (thin red curve), and additionally smoothed by taking the lagged average of the latest four forecasts at each start date (thick red curve, 676 members, see Methods). Forecast uncertainty (red shading, 5-95% confidence interval) is obtained from the forecast ensemble mean error variance (see Methods). (c) As (a) but for AMV and lagged ensemble. (d) As (c) but for NAO-matched forecast (see Methods). (e, f) As (c, d) but for northern European rainfall. Values of anomaly correlation (ACC) of the forecast ensemble mean and of persisting the latest observed 8-year mean available before each start date, and the ratio of predictable components (RPC), are indicated. Indices are defined in Methods. Time-series are anomalies relative to the average of all year 2-9 hindcasts.



**Figure 3: Decadal predictions of the extreme NAO period (1986 to 1997).** Observed anomalies of (a) temperature, (b) precipitation and (c) mean sea level pressure. (d, e, f) As (a, b, c) but for raw lagged ensemble mean forecasts. (g, h, i) As (d, e, f) but standardised by the ensemble mean standard deviation. (j, k, l) As (d, e, f) but for NAO-matched forecasts. Averages are taken for boreal winter (December to March) for all year 2-9 forecasts verifying in the period 1986 to 1997 (i.e. start dates 1985 to 1989 inclusive), and converted to anomalies by removing the average over all hindcasts (i.e. start dates 1960 to 2005 inclusive). Units are standard deviations. The raw lagged ensemble (d, e, f) is divided by the observed standard deviation to show the signal relative to observed variability.