

Improved Arctic sea ice thickness projections using bias corrected CMIP5 simulations

N. Melia¹, K. Haines² and E. Hawkins³

[1] Department of Meteorology, University of Reading, Reading, United Kingdom}

[2] National Centre for Earth Observation, Department of Meteorology, University of Reading, Reading, United Kingdom}

[3] NCAS-Climate, Department of Meteorology, University of Reading, Reading, United Kingdom}

Correspondence to: N. Melia (n.melia@pgr.reading.ac.uk)

Abstract

Projections of Arctic sea ice thickness (SIT) have the potential to inform stakeholders about accessibility to the region, but are currently rather uncertain. The latest suite of CMIP5 Global Climate Models (GCMs) produce a wide range of simulated SIT in the historical period (1979 – 2014) and exhibit various biases when compared with the Pan-Arctic Ice Ocean Modelling and Assimilation System (PIOMAS) sea ice reanalysis. We present a new method to constrain such GCM simulations of SIT via a statistical bias correction technique. The bias correction successfully constrains the spatial SIT distribution and temporal variability in the CMIP5 projections whilst retaining the climatic fluctuations from individual ensemble members. The bias correction acts to reduce the spread in projections of SIT and reveals the significant contributions of climate internal variability in the first half of the century and of scenario uncertainty from mid-century onwards. The projected date of ice-free conditions in the Arctic under the RCP8.5 high emission scenario occurs in the 2050s, which is a decade earlier than without the bias correction, with potentially significant implications for stakeholders in the Arctic such as the shipping industry. The bias correction methodology developed could be similarly applied to other variables to reduce spread in climate projections more generally.

1 **1 Introduction**

2 Global Climate Models (GCMs) are the primary tool for making climate predictions on
3 seasonal to decadal time scales, and climate projections over the next century (Flato et al.,
4 2013). In a warming climate, changes to sea ice thickness (SIT) are expected to lead to
5 significant implications for polar regions and beyond. A reduction in SIT will likely open up
6 the Arctic Ocean to economic diversification including new marine shipping routes (Smith
7 and Stephenson, 2013) and extraction of natural resources, as well as changes to the Arctic
8 ecosystem and potential links to mid-latitude weather (Francis and Vavrus, 2012). Many of
9 these economic opportunities may rely on SIT evolution, but current projections have
10 considerable uncertainty. SIT is also much more informative than sea ice concentration (SIC),
11 especially in the central Arctic, where future thinning can occur without major changes in the
12 local SIC.

13 The GCMs from the Coupled Model Intercomparison Project, phase 5 (CMIP5) (Taylor et al.,
14 2012) exhibit a large range in sea ice volume (SIV), spatial SIT distribution, and temporal SIT
15 variability under present day forcing conditions (e.g. Blanchard-Wrigglesworth and Bitz
16 (2014)). For September sea ice extent, Swart et al. (2015) showed the uncertainty in CMIP5
17 projections over the next few decades is dominated by these differences between models,
18 termed model uncertainty by Hawkins and Sutton (2009, 2011). Uncertainty in climate
19 projections arises from three distinct sources: (1) model uncertainty, (2) internal variability,
20 and (3) scenario uncertainty, as discussed by Hawkins and Sutton (2009, 2011) for
21 temperature and precipitation respectively. In contrast to projections of temperature where the
22 anomalies are often used, the absolute value of SIT is important – for example, ships have
23 critical SIT thresholds above which their use is not possible (Stephenson et al., 2013).

24 Bias correction (BC) of GCM simulations has the potential to reduce the differences between
25 models and hence potentially increase confidence in near term climate projections. The
26 importance of BC in impact based climate change studies was described in a special report of
27 the IPCC (Seneviratne et al., 2012), but BC has not previously been applied to projections of
28 SIT; this manuscript is novel in that it recalibrates SIT, and does it locally. There are many
29 different types of proposed BC techniques, (e.g. Boe et al. (2009); Christensen et al. (2008);
30 Ho et al. (2011); Mahlstein and Knutti (2012); Vrac and Friederichs (2014); Watanabe et al.
31 (2012), and references therein), which have mainly been applied to temperature and
32 precipitation. However, these existing methods need refining for sea ice as SIT is a

particularly challenging variable. This is due to its positive semi-definite nature, and the spatial and temporal occurrence of zeros, in observations and projections of SIT.

This study addresses the development of a new BC technique that constrains both the mean and variance of SIT in GCMs to an estimate of the observed statistics. It is important to correct the mean as this corrects the spatial SIT distribution. Variability in SIT also has a significant impact on the simulated range of regional ice-free dates, something of great interest to stakeholders, and the CMIP5 GCMs exhibit a wide range in their SIT variability. The study also uses multiple ensemble members from the same model when performing the BC, something that is often not utilised in other studies. This is important as it enables an assessment of the role of internal variability in future projections to be made. The techniques described in this paper are not limited to SIT, and would work for many climate variables. The exact implementation used in this study should also be calibrated to the user's needs based on factors such as the length of reliable observations and number of ensemble members.

In this paper we use the Pan-Arctic Ice Ocean Modelling and Assimilation System (PIOMAS) (Zhang and Rothrock, 2003) as a reanalysis based estimate of recent SIT, along with climate projections from a subset of six GCMs from the CMIP5 archive (Sect. 2). We first test the performance of increasingly complex BC approaches in a toy model environment (Sect. 3) and then apply our favoured method to the subset of CMIP5 GCMs in Sect. 4. We test the BC method by splitting the historical PIOMAS data, and then explore how the range in SIT projections is reduced using these techniques (Sect. 4) and summarise and discuss the results in Sect. 5.

2 Climate simulations and observations

2.1 PIOMAS

To represent observed SIT, we use estimates from the PIOMAS reanalysis. PIOMAS is a coupled ice-ocean model that is forced with the National Centers for Environmental Prediction (NCEP) atmospheric reanalysis, and assimilates satellite observed sea ice concentration (Lindsay and Zhang, 2006) and sea surface temperature (Schweiger et al., 2011). It does not however assimilate sea ice thickness (SIT), although this has been

attempted using the NASA Operation IceBridge and SIZONet campaigns of 2012 (Lindsay et al., 2012).

As a reanalysis, PIOMAS is constrained by the quality of the assimilated observations, Lindsay et al. (2014) forces PIOMAS with four different atmospheric reanalysis products producing differing results. Schweiger et al. (2011) found biases in PIOMAS of 0.26 m in autumn and 0.1 m in spring when compared with ICESat (Zwally et al., 2002) although the spring bias is within the range of uncertainties found by Zygmuntowska et al. (2014). Larger differences are found in the areas of thickest ice, north of Greenland and the Canadian Archipelago, with ICESat retrievals around 0.7 m larger than PIOMAS. However in this region PIOMAS agrees better with in situ data (Schweiger et al., 2011). Zygmuntowska et al. (2014) suggests that this discrepancy is due to the choice of sea ice density in ICESat, and they support this explanation by finding lower discrepancies between PIOMAS and CryoSat-2 (Laxon et al., 2013) which utilises an alternative sea ice density value. Stroeve et al. (2014), in a comprehensive study of SIT across CMIP5 and observations, find that the spatial correlations in thickness between CMIP5 models and PIOMAS are generally higher than those between CMIP5 models and ICESat. It should be noted that these results will be sensitive to the dataset chosen to represent observed SIT.

We choose PIOMAS to represent estimates of SIT as satellite observations are limited in their spatial and temporal range. For example, data from ICESat are only available between October and March 2003 – 2008 (Kwok et al., 2009). More recently Cryosat-2 has started producing real-time SIT datasets but only for the non-summer months (Tilling et al., 2015). This is also not ideal as it is the summer and autumn months when the ice is thinnest that are most relevant for potential economic activity. The spatial consistency, temporal length and completeness of the data are important considerations when computing climatological means and variances as the longest time series possible is needed to validate the statistics. It is for this reason primarily that PIOMAS has been chosen to represent observations in this study. Several studies (e.g. Laxon et al. (2013), Schweiger et al. (2011), Lindsay and Zhang (2006), and Stroeve et al. (2014)) have compared PIOMAS to satellite and in situ observations and models and find it a suitable estimate of observed SIT. PIOMAS is also deemed realistic enough to initialise numerical models for seasonal forecasts e.g., the Sea Ice Outlook (Blanchard-Wrigglesworth and Bitz, 2014) where the accuracy of the initial conditions is vital.

Figure 1 shows the mean September SIT and temporal standard deviation (SD) after linear detrending for PIOMAS over the satellite era (1979 – 2014). In the heart of the Canadian archipelago, PIOMAS ice thickness is up to 1.5 m, which is reasonable when compared to Haas and Howell (2015) who measured ice along the Northwest Passage in May 2011 and April 2015 using airborne electromagnetic induction soundings, and to Tilling et al. (2015) using Cryosat-2 for October and November 2010 – 2014. North of Greenland SIT exceeds 3.5 m, which is again comparable to Cryosat-2 for October and November 2010 – 2014 and is between zero and one meter along the north Russian coast. The SIT is most variable around the edge of the ice pack and especially near land. An effective BC should ensure that the simulations replicate these patterns of mean SIT and SD over this recent period.

2.2 Global climate models

This paper utilises a subset of six GCMs from CMIP5. Since a large part of this work assesses SIT variability, it is necessary for each GCM to have multiple ensemble simulations in the historical period and for each of the representative concentration pathways (RCPs) 2.6, 4.5 and 8.5 for future scenarios (Van Vuuren et al., 2011). In addition, the GCM mean spring thickness must fall within the 10th and 90th percentile of PIOMAS (Stroeve et al., 2014), have a reasonable spatial resolution, and a somewhat resolved Canadian archipelago. A consistent spatial distribution of land is needed for realistic and spatially complete multi-model means. The six GCMs that comprise this CMIP5 subset are listed in Table 1.

For the CMIP5 subset the historical simulations are used for the period 1979 – 2005. In most of the analysis for the period post-2005 the RCP8.5 scenario is used, which ramps up the amount of greenhouse gases to have a cumulative effect of increasing the direct radiative forcing by 8.5 Wm⁻² (approximately 1370 ppm CO₂ equivalent) by 2100 (Van Vuuren et al., 2011). The impact of other scenarios is compared later in the analysis. Figure 2 shows the 1979 – 2014 ensemble-mean September SIT for the CMIP5 subset, highlighting the considerable differences between the model simulations, and indicating that model bias is likely to be the dominant uncertainty in near-term projections.

The aim of the SIT BC outlined in this paper is to correct the mean and variance in the CMIP5 subset shown in Fig. 2 to the PIOMAS statistics. Although this should improve short-term predictions, a caveat to this approach is that PIOMAS only yields one realisation of the past (see Lindsay et al. (2014) for discussion of PIOMAS forced with alternative atmospheric

1 forcings). We have to assume that the relatively short period over which we have observations
2 (36 years) captures a representative sample of the behaviour we expect from the climate
3 system. In the short term, this is probably a reasonable assumption, as the GCMs will not
4 have evolved far from their corrected state of the recent past; this assumption is explored
5 further in Sect. 4.

7 **3 Bias correction methodology**

8 Bias correction methods effectively aim to reduce model uncertainty by constraining GCMs
9 to observations. There are two components to model uncertainty: the overall mean difference
10 (or bias), and differences in the amplitude of response to specified forcings. We have
11 deliberately chosen not to try and correct the simulated ice loss trend to that which PIOMAS
12 depicts. Our reasoning is to keep this as prescribed by the different GCMs because the
13 response of the SIT to future warming is unknown, likely non-linear, and the GCMs are
14 designed to give an estimate of this. It is also doubtful how well the forced current trend can
15 be determined from 36 years of data given the high noise to signal ratio for trends, especially
16 on grid point scales. It is also uncertain how much of the recent ice loss seen in the
17 observations can be attributed to changes in external forcing as opposed to internal variability,
18 although previous studies have attempted this including: Kay et al. (2011), Day et al. (2012),
19 Notz and Marotzke (2012), Stroeve et al. (2012), Notz (2015), Swart et al. (2015) and Zhang
20 (2015). We are also cautious of over fitting; applying a trend correction would potentially
21 result in an over-confident projection.

22 To test the performance of different possible BC methods a ‘toy model’ was used as proxy
23 ensemble timeseries (representing SIT at a single grid point for the same month each year for
24 the period 1979 – 2100). The timeseries are shown in Fig. 3a for a high mean - high variance
25 model (blue) and a low mean - low variance model (red), where the black line shows the
26 “truth” observations with one realisation over the historical period only. The time series were
27 all produced using a first order auto-regressive (with an AR(1) parameter of 0.3 chosen to be
28 representative of CMIP5 SIT auto-correlation) model imposed on a declining linear trend with
29 negative numbers reset to zero. Each model has five separate model ensemble members (thin
30 coloured lines) and the thick lines representing the ensemble means. The statistics in all the
31 legends are calculated over the observation window (1979 – 2014). ‘Ice-free’ in Fig. 3 is here
32 defined as the first occurrence of an ensemble member below 0.15 m. Shown is the ice-free

ensemble range, i.e. the year of the first ensemble member to be ice-free to the last ensemble member to be ice-free. A successful BC method should transform the individual ensemble members (thin red and blue lines) to match the mean and variance of the observations (black line), producing matched statistics. We test various approaches for such a bias correction. The mathematical notation for the following equations is in Table 2.

3.1 Additive correction

A basic additive correction, which has previously been used for temperature projections, is shown in Fig. 3b. This approach simply corrects the time-mean by subtracting the difference between the historical model ensemble-mean time-mean, $\langle \overline{M_h} \rangle$, and observation time mean, $\overline{O_h}$, from each of the model ensemble members, M .

$$\text{Additive corrected thickness} = M - (\langle \overline{M_h} \rangle - \overline{O_h}) \quad (1)$$

However, as the low ice model is adjusted up by the addition of a constant, it equilibrates at a positive value in the future rather than zero. Likewise the high ice model equilibrates at negative values. Neither of these properties are sensible.

This study makes use of multiple ensemble members from the same model, raising the question of how to treat ensemble member statistics when calculating a particular GCM's bias. For calculating the mean SIT, each GCM's ensemble mean is used because it is the GCM's mean bias that we wish to correct. This is important because a particular ensemble member's deviation from the ensemble mean is retained; it allows an individual ensemble member's time mean to be different to the observations over the historical period, but not the ensemble mean. The treatment of ensemble members for the SD calculation is described in section 3.4.

3.2 Multiplicative correction

If a multiplicative correction is used (Fig. 3c), where the ratio of the observed time mean and model ensemble-mean time-mean, $\overline{O_h} / \langle \overline{M_h} \rangle$, is multiplied as a factor to the model ensemble members, M , then the corrected thickness is:

$$\text{Multiplicative corrected thickness} = M \frac{\overline{O_h}}{\langle \overline{M_h} \rangle} \quad (2)$$

Multiplicative methods effectively preserve the future zero ice year, which is potentially an important value for a wide range of stakeholders. However, when applied as above this approach has the undesired effect of distorting the variances by the same factor as the mean correction, as visible in Fig. 3c.

3.3 Mean multiplicative correction

To avoid altering the variances, the mean multiplicative correction can be introduced (Fig. 3d), where the multiplicative mean correction, $\overline{O_h} / \langle \overline{M_h} \rangle$, is applied only to the 11-year-centred running-mean ensemble-mean, $\langle \tilde{M} \rangle$. This corrects the model mean evolution without corrupting the sub-decadal variance as $\langle \tilde{M} \rangle$ is smoothed. The model anomalies for each ensemble member, $M - \langle \tilde{M} \rangle$, are then added back to the corrected mean evolution:

$$\text{Mean multiplicative corrected thickness} = (M - \langle \tilde{M} \rangle) + \langle \tilde{M} \rangle \frac{\overline{O_h}}{\langle \overline{M_h} \rangle} \quad (3)$$

This works to correct the mean SIT and does not suffer from any peculiarities of the previous two methods. The model variance now remains unchanged but the approach opens up the possibility of correcting the variance towards that observed in the historical period. Note that by using the ensemble mean, $\langle \overline{M_h} \rangle$, for all these corrections we ensure that each ensemble member is corrected in the same way, thus preserving certain ensemble properties into the future.

3.4 Mean and variance correction

The GCMs from CMIP5 show a large range in SIT variance, and the magnitude of these variations is a significant factor determining when regions of the Arctic may first become accessible (when one ensemble member may first become ice-free). Therefore a variance correction is incorporated into Eq. (3) by taking the ratio of the temporal standard deviation of the detrended observations, $\sigma_{\overline{O_h}}$, to the square root of the ensemble mean of the variance of the detrended model ensembles, $\langle \sigma_{\tilde{M_h}} \rangle$ (detrended mean ensemble SD), over the historical period. The detrending in the models is calculated using each model's ensemble mean linear trend. This has some similarities to the approach of Ho et al. (2011) in application to temperature projections for Europe. Also see Appendix A for some further discussion of the choices made.

To incorporate the variance correction, the mean multiplicative correction (Eq. (3)) is first detrended, the variance correction applied, and the trend re-applied. This creates the Mean And VaRIance Correction (MAVRIC), shown in Eq. (4):

$$\text{MAVRIC} = (M - \langle \tilde{M} \rangle) \frac{\sigma_{\tilde{O}_h}}{\langle \sigma_{\tilde{M}_h} \rangle} + \langle \tilde{M} \rangle \frac{\overline{O_h}}{\langle \overline{M_h} \rangle} \quad (4)$$

Fig. 3e shows the MAVRIC does a near perfect job of correcting both the mean and variance to the observed statistics while still retaining the individual ensemble members' own climate fluctuations, but fractionally scaled by the variance ratio.

Comparing the ensemble range in projected ice-free date between the correction methods it is apparent that although the shapes of time-series have qualitatively changed this does not always result in a different range in projected ice-free date. For example on comparing the high mean – high variance GCM (blue) between (a) to (c) and (b) to (d); this is partly coincidence and partly due to how the four correction methods shown manipulate the time series. The MAVRIC method (e) results in a unique set of ice-free dates. This is an important attribute that the MAVRIC method displays, as the ice-free date is of vital importance to stakeholders in the Arctic and more basic methods of bias correction fail to appropriately adjust this parameter.

4 Bias corrected sea ice thickness projections

Figure 3e illustrates that the MAVRIC successfully corrects the mean and variance in a toy model environment. Before proceeding to investigate the impact of the MAVRIC on SIT projections it is prudent to test whether the MAVRIC can improve GCM performance by validating with PIOMAS. We use CSIRO-Mk3.6.0 (CSIRO) as the GCM to test. The ice in CSIRO generally has too much areal coverage and too little variability and is a CMIP5 outlier model with regards to SIT (Stroeve et al., 2014). However, CSIRO benefits from having 10 ensemble members, increasing the robustness of the statistics. For these two reasons, it is considered a thorough test of the MAVRIC's performance within a real GCM.

The test uses a data denial method where we train the MAVRIC on a subset of PIOMAS observations, 1979 – 1999, termed the calibration window. From this we examine how the MAVRIC predicts the observations for 2000 – 2014, termed the validation window. A limitation with this method is the length of observations: the period over which the MAVRIC calibration takes place must be long enough to capture a robust measure of the observed

1 statistics. The validation period must also be long enough to be able to draw robust
2 conclusions. It is not clear whether either the 21 year calibration or the 15 year validation
3 windows are long enough for robust method calibration and results verification, but we are
4 limited by the data available. An additional limitation to this method is that the calibration and
5 validation periods are very close to each other.

6 Figure 4 shows the performance of the MAVRIC at three grid points for September. The raw
7 CSIRO ensembles (grey) are bias corrected via the MAVRIC using the PIOMAS observations
8 (black) over the calibration window, producing the MAVRIC corrected ensembles (green) for
9 the validation window. If the MAVRIC can produce plausible predictions, the characteristics
10 of PIOMAS should be indistinguishable from individual corrected ensemble members in the
11 validation window. It is clear from the validation beanplots (right), that the distribution from
12 the corrected ensembles resembles PIOMAS much more closely than the raw distribution, e.g.
13 non-zero probability of zero ice. We do not expect the distribution from PIOMAS to match
14 the corrected distribution perfectly as PIOMAS only has one realisation (15 data points) while
15 CSIRO has 10 realisations. We can tentatively accept that this test demonstrates the validity
16 of the MAVRIC approach.

17 In the following sections the MAVRIC is applied to the CMIP5 subset of six GCMs used in
18 this study (Table 1). PIOMAS estimates of Arctic SIT are available from 1979 – 2014. This
19 36 year window is the period over which statistics are calculated in the observations, and in
20 the CMIP5 subset (using historical runs for 1979 – 2005 and RCP8.5 for 2006 – 2014). Each
21 model, month, and grid point has its own specific correction which is applied to all years
22 (1979 – 2100). However, separate ensemble members from the same GCM are treated with
23 the same correction, as we wish to correct the model bias and retain the ensemble spread.
24 Results are shown for September, initially only for CSIRO and later for all six models
25 combined to form the ‘CMIP5 subset’ used for this study.

26 **4.1 Temporal perspective example**

27 Figure 5 shows the impact of the MAVRIC in September in CSIRO at the same three grid
28 points as Fig. 4 but for the entire calibration window (1979 – 2014). The East Siberian Sea in
29 CSIRO has about double the SIT and half the SD of PIOMAS (Fig. 5a). The correction
30 therefore reduces the mean SIT whilst increasing the variance. This brings forward the range
31 of first year ice-free conditions (the first occurrence in each ensemble member of a SIT below

0.15 m) from after 2100 to 1981 – 2032. Similarly in the Beaufort Sea (Fig. 5b) the SD needs to be almost tripled, and the correction results in the first ice-free year coming over 100 years earlier. In the Fram Strait (Fig. 5c) CSIRO and PIOMAS have similar SIT requiring only a small mean adjustment, however CSIRO requires a big increase in variance. The MAVRIC moves the first possible ice-free date about 30 years earlier and increases the ensemble range from 32 to 63 years. It is worth noting that the dominant cause of this shift to earlier ice-free date at this location is due to the variance correction term in the MAVRIC rather than the mean correction term. This highlights the importance of correcting the variance in addition to the mean. Figure 5 demonstrates that the MAVRIC can lead to simulations that look significantly more like reality in the historical period and have an impact on regional ice-free projections.

4.2 Historical spatial perspective

In addition to examining the MAVRIC in a temporal sense, it is important to evaluate the results spatially to see where the MAVRIC is having the most effect and if it works at all locations. Figures 2 and 6 show that the mean September SIT distribution is very different in HadGEM2-ES and CSIRO. After the MAVRIC has been applied, the mean SIT fields are almost identical for the historical period (Fig. 6). It is important to note there are still differences when considering individual years and ensemble members i.e. the year-to-year variability and ensemble spread is preserved (although adjusted by the MAVRIC).

Figure 6 also shows the SD before and after the MAVRIC. The SD shown is the detrended mean ensemble SD as before. CSIRO has too low variability in the majority of locations although correctly places the maximum SD near the edges of the ice pack similarly to PIOMAS. HadGEM2-ES exhibits about the same magnitude of variability as the observations but the variability is too high in the centre of the ice pack and too low at the edges. After the correction the SD fields in both GCMs now look more similar to each other with the highest variability located at the edge of the ice pack and at coastal locations. They are now also both similar to the estimate from PIOMAS (Fig. 1).

4.3 CMIP5 subset multi-model sea ice thickness projections

The bias corrected SIT from each GCM can be brought together to form the multi-model mean CMIP5 subset, computed using three ensemble members (the maximum available

across all models) from each of the six GCMs for the historical and future decadal periods (Fig. 7). It is remarkable how the raw multi-model mean product for the historical period is not too different from PIOMAS in Fig 1, showing that the location and magnitude of model biases cancel out to a considerable degree, at least with this subset of models. Given this result it is not so surprising that the raw and corrected fields are fairly similar for the future projections also.

Nevertheless, even in this multi-model multi-ensemble framework the MAVRIC is still making some discernible differences. These differences are most apparent in the Canadian archipelago and the Russian Arctic seas, where the correction leads to a reduction in SIT of approximately 1 m in both regions. Both the raw and bias corrected fields predict a SIT loss of about 0.25 m per decade.

The fact that the MAVRIC is still making a significant difference on the regional scale is critical, e.g. for ship route availability. Currently studies that assess the future opening of Arctic shipping routes, which critically depend on the absolute value of SIT, do not yet account for such factors and will need to be reassessed.

4.4 Sources of uncertainty in projections of sea ice thickness

The uncertainty in climate projections can be partitioned into three distinct sources: (1) model uncertainty: for the same radiative forcing different models simulate different mean distributions and temporal changes. (2) Internal variability: the natural fluctuations of the climate present with or without any anthropogenic induced changes to radiative forcing. (3) Scenario uncertainty: uncertainty in future radiative forcing resulting from unknown future emissions. Hawkins and Sutton (2009, 2011) assessed these sources of uncertainty in global and regional temperature and precipitation projections, and here we quantify the sources of uncertainty in SIT, utilising the CMIP5 subset multi-model ensemble. Crucially we use the absolute values of SIT rather than considering anomalies as is often done for other climate variables. The methodology for partitioning these sources of uncertainty is detailed in Appendix B. An additional source of uncertainty that we neglect here is the PIOMAS calibration uncertainty emerging from the choice of atmospheric reanalysis and model tuning. This could be assessed by sampling the different versions of the PIOMAS reanalysis described in Lindsay et al. (2014). They find the different versions are broadly similar and can be accounted for by appropriate tuning of the ice model component. This bias in PIOMAS

1 itself will introduce systematic biases to the MAVRIC projections. This bias is not a flaw in
2 MAVRIC however but a limitation intrinsic to the observational dataset one is correcting to.

3 The MAVRIC method outlined in this study acts to eliminate the model bias in the MAVRIC
4 calibration period (1979 – 2014). After this period the model uncertainty grows due to the
5 GCM's differing responses to changes in external forcing. The sources of uncertainty for SIT
6 for the decade 2015 – 2024, immediately following the MAVRIC calibration period, are
7 shown in Fig. 8. The total uncertainty in the corrected CMIP5 subset is strikingly lower than
8 in the raw CMIP5 subset. Closer analysis reveals that this is due to the substantial reduction in
9 model uncertainty owing to the MAVRIC. The other sources of uncertainty do not change as
10 much.

11 The temporal evolution of these sources of uncertainty is shown in Fig. 9a by taking the
12 median variance from each of the panels in Fig. 8 for this and other periods. There are three
13 competing factors for how the uncertainty will change with time. First, the SIT is decreasing,
14 and this will reduce the uncertainty as the range of values of which the SIT can occupy
15 shrinks. Second, the separate GCM's simulated SIT responses due to external forcing will
16 differ from each other, causing GCMs to drift apart over time. Thirdly, sea ice at the grid
17 point scale becomes more mobile and vulnerable to external factors as it thins. This will
18 increase variability, initially at least (Sou and Flato, 2009). All of these factors are involved in
19 the evolution of the uncertainties.

20 The raw CMIP5 subset exhibits a decrease in total uncertainty with time (dashed black in Fig.
21 9a). This is primarily due to the reduction in model uncertainty (dashed blue), likely because
22 the mean SIT is reducing. The corrected total uncertainty is lower than the raw uncertainty
23 until at least the end of the century. This means that the MAVRIC can reduce the model
24 spread (or bias) and so may potentially increase confidence in climate projections of SIT
25 throughout this period. The corrected model uncertainty increases for the first three decades,
26 as the models start from a similar state and subsequently diverge because of differing
27 responses to the changes in external forcing. Later the corrected model uncertainty reduces as
28 the mean SIT decreases towards zero.

29 The total uncertainty is the sum of model uncertainty, internal variability, and scenario
30 uncertainty (see Appendix B for more details). The other panels in Fig. 9 illustrate the relative
31 importance of these sources of uncertainty in terms of the percentage total variance explained,
32 for the raw data, and after the MAVRIC.

Fig. 9b illustrates that in the raw projections, model uncertainty remains the dominant (> 50 %) source of uncertainty until at least 2100, whereas it only becomes dominant for a few decades mid-century after the MAVRIC (Fig. 9c). The absolute magnitude of internal variability, and its contribution to the total uncertainty, decreases with time because SIT also decreases with time. In the corrected projections, the internal variability is the major contributor to the total uncertainty for the first 25 years, compared to a maximum contribution of only 26 % in the raw projections. This highlights the importance of correcting the variance to realistic magnitudes and also the key role of natural variations in predicting the near future evolution of sea ice. The scenario uncertainty accounts for less than 10 % of the total uncertainty for the first 50+ years. Additional analysis metrics on the improvement the MAVRIC method affords can be found in Appendix C.

Although we have demonstrated here that the MAVRIC method reduces the model uncertainty as seen by the reduction in spread of projected SIT with our selection of GCMs, we acknowledge that this may not necessarily correspond to a reduction in uncertainty in the real world.

4.5 Reduced spread in timing of ice-free conditions

By reducing the model spread the range of possible outcomes has been reduced, this potentially leads to greater confidence in SIT projections. Figure 10 shows the raw and corrected CMIP5 subset SIV* projections until 2100 using the 18 multi-model ensemble members in each scenario as before (* calculated here does not consider SIC as it is not bias corrected). To find a representative SIC for the SIV* calculation we use the September SIC in CCSM4 RCP8.5 and find a mean (of the non-zero grid cells) SIC of approximately 50% for 2006-2100.

The thick coloured lines are the multi-model scenario mean and the coloured regions represent the 16 – 84 percentiles (equivalent to 1σ around the mean of a Gaussian distribution) of the ensemble members. To account for the large range in SIT at any particular time in the CMIP5 subset, we use a method similar to that of Massonnet et al. (2012) to calculate first ice-free conditions. We postulate that SIV for ice-free conditions is $1 \times 10^3 \text{ km}^3$, which is in agreement with previous studies calculating first ice-free dates (e.g. Massonnet et al. (2012) and Overland and Wang (2013)), and is equivalent to one meter thick ice for an ice extent of 10^6 km^2 .

The MAVRIC reduces the total SIV, but the relative magnitude of this reduction decreases as SIV declines. The 16 – 84 % range has also been vastly reduced, particularly for the near future. For example, in 2025 the MAVRIC has reduced the 16 – 84 % range from $6 \times 10^3 \text{ km}^3$ to $2.5 \times 10^3 \text{ km}^3$. It is this reduction in the plausible range of SIV that leads to potential increased confidence in projections of SIT and SIV. To assess when the Arctic will first display ice-free conditions, we focus on RCP8.5, the most realistic scenario from the last 10 years (Fuss et al., 2014). The cumulative number of ensemble members having satisfied the ice-free criterion as a function of time is shown in Fig. 10c. If the range in this parameter has reduced, this will be shown by the gradient of the line increasing after MAVRIC, and this is clearly seen. Figure 10d further illustrates the spread reduction with boxplots, where the line represents the median (9th) ensemble member to go ice-free. This occurs in 2052 with the MAVRIC, nine years earlier than before. The box represents 16 – 84 % of the ensemble members, this range has been reduced by about 20 years; dates after 2085 can now be eliminated.

Corrected results from the other emission scenarios show similar features but with later ice-free dates, as expected for lower emissions, and some ensemble members fail to go ice-free by 2100. For RCP4.5 the MAVRIC makes a profound difference with the median ice-free date occurring 35 years earlier in 2060. For RCP2.6 there is spread reduction mid-century but the CMIP5 subset before and after the MAVRIC are in good agreement by the end of the century, with projected ice-free dates around 2090.

5 Summary and discussion

5.1 Summary

This study has developed a bias correction methodology for simulations of sea ice thickness (SIT). By constraining CMIP5 simulations with the PIOMAS reanalysis we have demonstrated that:

- GCMs simulate a wide range of SIT in the historical period and exhibit various spatial and temporal biases when compared with the PIOMAS reanalysis. This model uncertainty is the dominant source of uncertainty in CMIP5 future climate projections of SIT.
- The Mean And VaRIance Correction (MAVRIC) technique outlined in this paper significantly reduces the total uncertainty in future projections of SIT out to 2100 by

reducing model uncertainty. Correcting both mean and variance of models is found to be critical for improving the robustness of the projections.

- The MAVRIC results in internal variability being the dominant source of uncertainty until 2022, and model uncertainty is dominant thereafter. From mid-century onwards, scenario uncertainty becomes increasingly important and as influential as model uncertainty by 2100.
- The MAVRIC results in projected September ice-free conditions in the Arctic under RCP8.5 occurring up to 10 years earlier (2050s) than without the correction, and with a much narrower range, e.g. excluding post 2085 dates.

5.2 Discussion

Without the MAVRIC, the true magnitude of the internal variability and scenario uncertainty in projections of SIT is concealed by the dominant model uncertainty. This demonstrates that time invested in running many ensemble members to sample internal variability in SIT may be more beneficial than running many future emission scenarios for near term projections. These findings implicate that there is room for improvement in GCMs at least for 50 year projections where the scenario differences are negligible. However, for projections at the end of the century, the scenarios become more important.

The MAVRIC bias correction technique developed in this study results in a significant improvement in model simulations of SIT with respect to observations. In future projections, the MAVRIC results in a substantial reduction in the range of SIT, potentially leading to increased confidence in climate projections. As absolute values of SIT are utilised, this reduction in spread potentially has important implications for stakeholder sectors operating in Arctic waters such as shipping. The application of the bias correction results in a 60% reduction in the likely range (16 – 84 percentiles) of sea ice volume in September 2025.

There are a number of caveats to these findings. No attempt is made to constrain the trend in the GCMs. This would be difficult because of the short time scale over which observations are available, raising serious questions about the robustness of calculated historical trends. However future studies could consider this further and assess the feasibility of a trend correction to GCMs. In addition, it is important to recognise that PIOMAS, used here as observations, will also have errors. It would be possible to reduce the multiplicative

1 weightings in Eq. (4) to reflect some uncertainty in the historical data. Other temporally and
2 spatially complete sea ice reanalyses could also be used in future to address this issue.

3 The simulations tend to show an increase in variance as the sea ice thins, before subsequently
4 declining as the thickness approaches zero (Goosse et al., 2009). Blanchard-Wrigglesworth
5 and Bitz (2014) assessed the relationship of this mean state dependant variance in 19 GCMs,
6 including five of the six used in this study, in addition to PIOMAS. They find a relationship
7 between mean thickness variability and mean thickness in models, i.e. models with thicker
8 SIT depict more variable SIT. In the 19 GCMs assessed, PIOMAS sits on the trend line for
9 the correlation between mean thickness variability and mean thickness. However, in the
10 developed MAVRIC, the change in variance is decoupled from the applied change to the
11 mean state. This aspect could be further developed, but only by making additional
12 assumptions about future changes in SIT variability. Studies should make use of the MAVRIC
13 in assessing the impact on potential stakeholders sensitive to SIT and a paper utilising the
14 MAVRIC to investigate the opening of the Arctic sea routes is in preparation. We also make
15 the bias corrected SIT fields (Melia, 2015), freely available online for further investigations at
16 <http://dx.doi.org/10.17864/1947.9>.

Appendix A Supplementary MAVRIC methodology details

For model biases to be calculated a common grid needed to be used, hence all MAVRIC calculations took place on the CMIP5 model's native grid. This means that PIOMAS was converted to the CMIP5 model grid for each GCM's bias calculations. This choice was made as it only involves interpolating one of the two fields each time and generally it is PIOMAS that has the higher resolution. The BC shown in Eq. (4) contains two terms for the representation of the variance in both observations $\sigma_{\widehat{\sigma}_h}$ and models $\langle \sigma_{\widehat{M}_h} \rangle$. Over the 36 year period of observations the magnitude of the ice loss trend can be significant. To accurately calculate variances this externally forced trend should first be removed to leave the variance due to internal variability. Here a choice needs to be made about how best to remove the externally forced trend. For the PIOMAS observations we choose to linearly detrend the monthly data. A smoothed detrending was considered, however this might remove longer time scale variability which is undesirable. Using similar reasoning it is possible that the linear detrending is removing some variability on the multi-decadal timescale. This is assumed to be significantly less than variability on smaller timescales, and much of the trend is attributed to be externally forced over the 36 years, hence should not be included as internal variability. The performance of a smoothed detrend was tested in a theoretical framework and resulted in a 10 % loss of accuracy in the standard deviation correction due to describing variance as trend.

The calculation of variance in the models is more complicated due to the fact that there is more than one realisation. It is obvious that the required variance should be calculated from the individual ensemble members rather than the ensemble mean. The variance should be calculated in each ensemble member and then the mean taken. There is another choice to make, i.e. whether each ensemble member should be detrended with its own trend, or should the ensemble mean trend be used? We propose that the ensemble mean trend should be used as this is the models response to the changes in forcings. The model detrended ensemble mean standard deviation, $\langle \sigma_{\widehat{M}_h} \rangle$, was calculated by calculating the detrended ensemble variances, then taking the square root of their mean.

The running mean for the future model correction term $\langle \widetilde{M} \rangle$ is calculated over an 11 year period of the ensemble mean, this window hence starts at 1975 for the historical calculations. The chosen period must be long enough to adequately smooth the time series, whilst still

- 1 being able to capture variations in the sea ice decline trend. This was also tested and found to
- 2 outperform a 21 year period.
- 3

Appendix B Partitioning sources of uncertainty

The sources of uncertainty in Sect. 4.4, Figs. 8 and 9 are calculated for each decadal period (2005 – 2014, 2015 – 2024, etc.) separately as follows. Three ensemble members from each of the six GCMs are utilised for three different emission scenarios (RCP2.6, 4.5, and 8.5). This results in each decade having $6(\text{GCMs}) \times 3(\text{ensemble members}) \times 3(\text{scenarios}) \times 10(\text{years}) = 540(\text{fields})$.

- The total uncertainty is the variance calculated across all 540 fields.
- The internal variability is calculated similarly to the total variability except instead of the absolute values the anomalies from the models' decadal-mean ensemble-mean for each scenario are used.
- To calculate the model uncertainty, each of the six models' decadal-mean ensemble-mean is calculated, resulting in six fields. The variance is then calculated across these six fields, and repeated for all three scenarios separately (to eliminate differential model dependent responses to the different emission scenarios). The model uncertainty is the square root of the mean of these three fields.
- The scenario uncertainty is calculated in a similar way. For each model, each of the three scenarios decadal-mean ensemble-means are calculated resulting in three (scenario-dependant) decadal-mean ensemble-means for each of the six models. The variance is then calculated through these three scenario mean fields for each of the six models, resulting in six fields of the variance in each model. The square root of the mean of the six models scenario uncertainty is the scenario uncertainty.

To create Fig. 8b and c it is assumed that the total variance (total uncertainty, T^2) is the sum of the variance due to model uncertainty (M^2), internal variability (I^2), and scenario uncertainty (S^2), formally:

$$T^2 = M^2 + I^2 + S^2 \quad (\text{B1})$$

We note that the variances calculated above do not always sum exactly in this way due to small interaction terms (approximately 10%) which we ignore.

1 **Appendix C Additional MAVRIC performance analysis**

2 To highlight whether the estimated uncertainties are reliable, we examine the errors in the
3 projections when considering one member as ‘truth’. As all ensemble members are
4 constrained by PIOMAS one individual ensemble member out of sample should fall with in
5 the distribution of the remaining ensemble members. This principle should hold true for all
6 ensemble members out of sample in turn.

7 The root mean square error (RMSE) is calculated using the Eq. (C1):

$$RMSE = \sqrt{\frac{1}{18} \sum_{n=1}^{18} (E_n - \overline{E_{15}})^2} \quad (C1)$$

8 where E_n is the ensemble member between 1 to 18, $\overline{E_{15}}$ is the mean of the 15 ensemble
9 members from the models of which E_n is not a member.

10 Figure C1 shows the advantage of the MAVRIC method in this out of sample RMSE test. A
11 decreasing RMSE means that the models are initially biased though are converging to a
12 common value (as we expect in this case as the models trend towards being ice-free). An
13 increasing RMSE means that the models are diverging as they have different ice loss trends.

14 Figure C1 shows the advantage of the MAVRIC method in this out of sample RMSE test. A
15 decreasing RMSE means that the models are initially biased though are converging to a
16 common value (as we expect in this case as the models trend towards being ice-free). An
17 increasing RMSE means that the models are diverging as they have different ice loss trends.

18 The MAVRIC ensemble trained on every individual ensemble member within MAVRIC
19 results in a RMSE of 0.1 m initially and up to a maximum RMSE of 0.5 m. The fact that the
20 Raw RMSE decreases (as opposed to increases) highlights that the models have biases. The
21 0.1 m in the MAVRIC RMSE indicates that initially the MAVRIC ensemble members differ
22 only in internal variability. The RMSE then grows due to differing ice loss trends which is
23 expected as no attempt to correct the trends in this study.

24 To find the dispersion of the MAVRIC multi-model ensemble we repeat this style of
25 experiment with the standard error (SE) metric, using Eq (C2):

$$SE = \frac{E_n - \overline{E_{15}}}{\sigma_{15}} \quad (C1)$$

1 where E_n is the ensemble member between 1 to 18, $\overline{E_{15}}$ is the mean of the 15 ensemble
2 members from the models of which E_n is not a member. σ_{15} is the standard deviation of the
3 15 ensemble members of which E_n is not a member. This is repeated for all 18 ensemble
4 members giving 18 SEs of how different each ensemble member is to the rest of the multi-
5 model ensemble set. The SD across these 18 SEs is the dispersion of the multi-model
6 ensemble. A perfectly dispersed ensemble set will have a dispersion of one. Numbers less
7 than one mean the ensemble set is under-dispersed and hence predictions/projections from
8 that set will be under-confident as the SD is too large. Values greater than one indicate that
9 the system is over-dispersive and hence over-confident.

10 The results of the dispersion calculation are shown in Fig. C2. The MAVRIC ensemble is
11 approximately 15 % - 30 % over-dispersed for lead times of up to 60 years. This means that
12 the ensemble is slightly over-confident and thus has slightly too little overall variance. The
13 rapid increase in dispersion from 60 years is solely due to the CSIRO GCM, specifically it's
14 comparatively slow ice loss trend. This was tested by repeating the dispersion experiment
15 omitting CSIRO (not shown). At this lead time many models are starting to be ice-free in
16 September while CSIRO retains ice. It is to the merit of MAVRIC that it is less over-
17 dispersed than the Raw output, hence more reliance can be placed on MAVRIC than the Raw
18 output as it's ensemble distribution is more representative.

Author contribution

N. M., K. H., and E. H. designed the methodology and experiments.

N.M. developed the code, and performed the experiments.

N. M., K. H., and E. H. wrote the manuscript.

Acknowledgements

We thank Dr Steffen Tietsche for the conversion of the PIOMAS data, Prof. Daniel Feltham and Prof. Ellie Highwood for comments on a pre-submission draft. We thank Referees Prof. Gregory Flato and Dr Francois Massonnet for their quick responses and thorough and constructive reviews. We thank Dr Robert Darby and the University of Reading Research Data Archive for facilitating the hosting of the MAVRIC data set. All statistical analyses and figures were accomplished using the R language and environment for statistical computing and graphics. For more information, see <http://www.r-project.org/>.

N. M. and E. H. are funded by the APPOSITE project (grant NE/I029447/1), funded by the UK Natural Environment Research Council as part of the Arctic Research Programme. E. H. is also funded by NERC Fellowship. K. H. is partly funded by the National Centre for Earth Observation NCEO.

References

- Blanchard-Wrigglesworth, E. and Bitz, C. M.: Characteristics of Arctic Sea-Ice Thickness Variability in GCMs, *J. Clim.*, 27, 8244-8258, doi: 10.1175/Jcli-D-14-00345.1, 2014.
- Boe, J., Hall, A., and Qu, X.: September sea-ice cover in the Arctic Ocean projected to vanish by 2100, *Nat. Geosci.*, 2, 341-343, doi: 10.1038/ngeo467, 2009.
- Christensen, J. H., Boberg, F., Christensen, O. B., and Lucas-Picher, P.: On the need for bias correction of regional climate change projections of temperature and precipitation, *Geophys. Res. Lett.*, 35, L20709, doi: 10.1029/2008gl035694, 2008.
- Day, J. J., Hargreaves, J. C., Annan, J. D., and Abe-Ouchi, A.: Sources of multi-decadal variability in Arctic sea ice extent, *Environ. Res. Lett.*, 7, 034011, doi: 10.1088/1748-9326/7/3/034011, 2012.
- Flato, G., Marotzke, J., Abiodun, B., Braconnot, P., Chou, S. C., Collins, W., Cox, P., Driouech, F., Emori, S., Eyring, V., Forest, C., Gleckler, P., Guilyardi, E., Jakob, C., Kattsov, V., Reason, C., and Rummukainen, M.: Evaluation of Climate Models. In: *Climate Change 2013: The Physical Science Basis. Contribution of Working Group I to the Fifth Assessment Report of the Intergovernmental Panel on Climate Change*, edited by: Stocker, T. F., Qin, D., Plattner, G.-K., Tignor, M., Allen, S. K., Boschung, J., Nauels, A., Xia, Y., Bex, V., and Midgley, P. M., Cambridge University Press, Cambridge, United Kingdom and New York, NY, USA, 741–866, 2013.
- Francis, J. A. and Vavrus, S. J.: Evidence linking Arctic amplification to extreme weather in mid-latitudes, *Geophys. Res. Lett.*, 39, L06801, doi: 10.1029/2012gl051000, 2012.
- Fuss, S., Canadell, J. G., Peters, G. P., Tavoni, M., Andrew, R. M., Ciais, P., Jackson, R. B., Jones, C. D., Kraxner, F., Nakicenovic, N., Le Quere, C., Raupach, M. R., Sharifi, A., Smith, P., and Yamagata, Y.: Betting on negative emissions, *Nat Clim Change*, 4, 850-853, doi: 10.1038/nclimate2392, 2014.
- Gent, P. R., Danabasoglu, G., Donner, L. J., Holland, M. M., Hunke, E. C., Jayne, S. R., Lawrence, D. M., Neale, R. B., Rasch, P. J., and Vertenstein, M.: The community climate system model version 4, *J. Clim.*, 24, 4973-4991, 2011.
- Goosse, H., Arzel, O., Bitz, C. M., de Montety, A., and Vancoppenolle, M.: Increased variability of the Arctic summer ice extent in a warmer climate, *Geophys. Res. Lett.*, 36, L23702, doi: 10.1029/2009gl040546, 2009.
- Haas, C. and Howell, S. E. L.: Ice thickness in the Northwest Passage, *Geophys. Res. Lett.*, 42, 7673-7680, doi: 10.1002/2015gl065704, 2015.
- Hawkins, E. and Sutton, R.: The Potential to Narrow Uncertainty in Regional Climate Predictions, *Bull. Am. Meteorol. Soc.*, 90, 1095-1107, doi: 10.1175/2009bams2607.1, 2009.
- Hawkins, E. and Sutton, R.: The potential to narrow uncertainty in projections of regional precipitation change, *Clim. Dyn.*, 37, 407-418, doi: 10.1007/s00382-010-0810-6, 2011.
- Ho, C. K., Stephenson, D. B., Collins, M., Ferro, C. A. T., and Brown, S. J.: Calibration Strategies: A Source of Additional Uncertainty in Climate Change Projections, *Bull. Am. Meteorol. Soc.*, 93, 21-26, doi: 10.1175/2011bams3110.1, 2011.

1 Jungclaus, J., Keenlyside, N., Botzet, M., Haak, H., Luo, J.-J., Latif, M., Marotzke, J.,
2 Mikolajewicz, U., and Roeckner, E.: Ocean circulation and tropical variability in the coupled
3 model ECHAM5/MPI-OM, *J. Clim.*, 19, 3952-3972, 2006.

4 Kay, J. E., Holland, M. M., and Jahn, A.: Inter-annual to multi-decadal Arctic sea ice extent
5 trends in a warming world, *Geophys. Res. Lett.*, 38, L15708, doi: 10.1029/2011gl048008,
6 2011.

7 Kwok, R., Cunningham, G. F., Wensnahan, M., Rigor, I., Zwally, H. J., and Yi, D.: Thinning
8 and volume loss of the Arctic Ocean sea ice cover: 2003-2008, *J. Geophys. Res. Oceans*, 114,
9 C07005, doi: 10.1029/2009jc005312, 2009.

10 Laxon, S. W., Giles, K. A., Ridout, A. L., Wingham, D. J., Willatt, R., Cullen, R., Kwok, R.,
11 Schweiger, A., Zhang, J., Haas, C., Hendricks, S., Krishfield, R., Kurtz, N., Farrell, S., and
12 Davidson, M.: CryoSat-2 estimates of Arctic sea ice thickness and volume, *Geophys. Res.*
13 *Lett.*, 40, 732-737, doi: 10.1002/Grl.50193, 2013.

14 Lindsay, R., W., Haas, C., Hendricks, S., Hunkeler, P., Kurtz, N., Paden, J., Panzer, B.,
15 Sonntag, J., Yungel, J., and Zhang, J.: Seasonal forecasts of Arctic sea ice initialized with
16 observations of ice thickness, *Geophys. Res. Lett.*, 39, L21502, doi: 10.1029/2012gl053576,
17 2012.

18 Lindsay, R., Wensnahan, M., Schweiger, A., and Zhang, J.: Evaluation of Seven Different
19 Atmospheric Reanalysis Products in the Arctic*, *J. Clim.*, 27, 2588-2606, doi: 10.1175/jcli-d-
20 13-00014.1, 2014.

21 Lindsay, R. W. and Zhang, J.: Assimilation of Ice Concentration in an Ice–Ocean Model, *J.*
22 *Atmos. Oceanic Technol.*, 23, L21502, doi: 10.1029/2012gl053576, 2006.

23 Mahlstein, I. and Knutti, R.: September Arctic sea ice predicted to disappear near 2°C global
24 warming above present, *J. Geophys. Res. Atmos.*, 117, D06104, doi: 10.1029/2011jd016709,
25 2012.

26 Massonnet, F., Fichefet, T., Goosse, H., Bitz, C. M., Philippon-Berthier, G., Holland, M. M.,
27 and Barriat, P. Y.: Constraining projections of summer Arctic sea ice, *The Cryosphere*, 6,
28 1383-1394, doi: 10.5194/tc-6-1383-2012, 2012.

29 Meehl, G. A., Washington, W. M., Arblaster, J. M., Hu, A., Teng, H., Kay, J. E., Gettelman,
30 A., Lawrence, D. M., Sanderson, B. M., and Strand, W. G.: Climate change projections in
31 CESM1 (CAM5) compared to CCSM4, *J. Clim.*, 26, 6287-6308, 2013.

32 Melia, N.: Improved Arctic sea ice thickness projections using bias corrected CMIP5
33 simulations., University of Reading, doi: <http://dx.doi.org/10.17864/1947.9>, 2015.

34 Notz, D.: How well must climate models agree with observations?, *Philosophical*
35 *Transactions of the Royal Society of London A: Mathematical, Physical and Engineering*
36 *Sciences*, 373, doi: 10.1098/rsta.2014.0164, 2015.

37 Notz, D. and Marotzke, J.: Observations reveal external driver for Arctic sea-ice retreat,
38 *Geophys. Res. Lett.*, 39, L08502, doi: 10.1029/2012gl051094, 2012.

39 Overland, J. E. and Wang, M.: When will the summer Arctic be nearly sea ice free?,
40 *Geophys. Res. Lett.*, 40, 2097-2101, doi: 10.1002/grl.50316, 2013.

41 Rotstayn, L., Jeffrey, S., Collier, M., Dravitzki, S., Hirst, A., Syktus, J., and Wong, K.:
42 Aerosol-and greenhouse gas-induced changes in summer rainfall and circulation in the

1 Australasian region: a study using single-forcing climate simulations, *Atmos. Chem. Phys.*, 12,
2 6377-6404, doi: 10.5194/acp-12-6377-2012, 2012.

3 Schweiger, A., Lindsay, R., Zhang, J., Steele, M., Stern, H., and Kwok, R.: Uncertainty in
4 modeled Arctic sea ice volume, *J. Geophys. Res. Oceans*, 116, doi: 10.1029/2011jc007084,
5 2011.

6 Seneviratne, S. I., Nicholls, N., Easterling, D., Goodess, C. M., Kanae, S., Kossin, J., Luo, Y.,
7 Marengo, J., McInnes, K., and Rahimi, M.: Changes in climate extremes and their impacts on
8 the natural physical environment. In: *Managing the risks of extreme events and disasters to
9 advance climate change adaptation*, edited by: Field, C. B., Barros, V., Stocker, T. F., Qin, D.,
10 Dokken, D. J., Ebi, K. L., Mastrandrea, M. D., Mach, K. J., Plattner, G. K., K., A. S., Tignor,
11 M., and M., M. P., A Special Report of Working Groups I and II of the Intergovernmental
12 Panel on Climate Change (IPCC), Cambridge University Press, Cambridge, UK, and New
13 York, NY, USA, 109-230, 2012.

14 Smith, L. C. and Stephenson, S. R.: New Trans-Arctic shipping routes navigable by
15 midcentury, *Proc. Natl. Acad. Sci. U.S.A.*, 110, E1191–E1195, doi:
16 10.1073/pnas.1214212110, 2013.

17 Sou, T. and Flato, G.: Sea Ice in the Canadian Arctic Archipelago: Modeling the Past (1950–
18 2004) and the Future (2041–60), *J. Clim.*, 22, 2181-2198, doi: 10.1175/2008jcli2335.1, 2009.

19 Stephenson, S., Smith, L., Brigham, L., and Agnew, J.: Projected 21st-century changes to
20 Arctic marine access, *Clim. Change*, 118, 885-899, doi: 10.1007/s10584-012-0685-0, 2013.

21 Stroeve, J., Barrett, A., Serreze, M., and Schweiger, A.: Using records from submarine,
22 aircraft and satellite to evaluate climate model simulations of Arctic sea ice thickness, *The
23 Cryosphere*, 8, 1839-1845, doi: 10.5194/tc-8-1839-2014, 2014.

24 Stroeve, J., Serreze, M., Holland, M., Kay, J., Malanik, J., and Barrett, A.: The Arctic's
25 rapidly shrinking sea ice cover: a research synthesis, *Clim. Change*, 110, 1005-1027, doi:
26 10.1007/s10584-011-0101-1, 2012.

27 Swart, N. C., Fyfe, J. C., Hawkins, E., Kay, J. E., and Jahn, A.: Influence of internal
28 variability on Arctic sea-ice trends, *Nat Clim Change*, 5, 86-89, doi: 10.1038/nclimate2483
29 2015.

30 Taylor, K. E., Stouffer, R. J., and Meehl, G. A.: An overview of CMIP5 and the experiment
31 design, *Bull. Am. Meteorol. Soc.*, 93, 485-498, doi: 10.1175/Bams-D-11-00094.1, 2012.

32 The HadGEM2 Development Team, Martin, G. M., Bellouin, N., Collins, W. J., Culverwell,
33 I. D., Halloran, P. R., Hardiman, S. C., Hinton, T. J., Jones, C. D., McDonald, R. E.,
34 McLaren, A. J., O'Connor, F. M., Roberts, M. J., Rodriguez, J. M., Woodward, S., Best, M. J.,
35 Books, M. E., Brown, A. R., Butchart, N., Dearden, C., Derbyshire, S. H., Dharssi, I.,
36 Doutriaux-Boucher, M., Edwards, J. M., Falloon, P. D., Gedney, N., Grey, L. J., Hewitt, H.
37 T., Hobson, M., Huddleston, M. R., Hughes, J., Ineson, S., Ingram, W. J., James, P. M., Johns,
38 T. C., Johnson, C. E., Jones, A., Jones, C. P., Joshi, M. M., Keen, A. B., Liddicoat, S., Lock,
39 A. P., Maidens, A. V., Manners, J. C., Milton, S. F., Rae, J. G. L., Ridley, J. K., Sellar, A.,
40 Senior, C. A., Totterdell, I. J., Verhoef, A., Vidale, P. L., and Wiltshire, A.: The HadGEM2
41 family of Met Office Unified Model climate configurations, *Geosci. Model Dev.*, 4, 723-757,
42 doi: 10.5194/gmd-4-723-2011, 2011.

- 1 Tilling, R. L., Ridout, A., Shepherd, A., and Wingham, D. J.: Increased Arctic sea ice volume
2 after anomalously low melting in 2013, *Nat. Geosci*, 8, 643-646, doi: 10.1038/ngeo2489,
3 2015.
- 4 Van Vuuren, D. P., Edmonds, J., Kainuma, M., Riahi, K., Thomson, A., Hibbard, K., Hurtt,
5 G. C., Kram, T., Krey, V., and Lamarque, J.-F.: The representative concentration pathways:
6 an overview, *Clim. Change*, 109, 5-31, doi: 10.1007/s10584-011-0148-z, 2011.
- 7 Vrac, M. and Friederichs, P.: Multivariate—Intervariable, Spatial, and Temporal—Bias
8 Correction, *J. Clim.*, 28, 218-237, doi: 10.1175/jcli-d-14-00059.1, 2014.
- 9 Watanabe, M., Suzuki, T., O'ishi, R., Komuro, Y., Watanabe, S., Emori, S., Takemura, T.,
10 Chikira, M., Ogura, T., and Sekiguchi, M.: Improved climate simulation by MIROC5: mean
11 states, variability, and climate sensitivity, *J. Clim.*, 23, 6312-6335, doi:
12 10.1175/2010jcli3679.1, 2010.
- 13 Watanabe, S., Kanae, S., Seto, S., Yeh, P. J. F., Hirabayashi, Y., and Oki, T.: Intercomparison
14 of bias-correction methods for monthly temperature and precipitation simulated by multiple
15 climate models, *J. Geophys. Res. Atmos.*, 117, doi: 10.1029/2012jd018192, 2012.
- 16 Zhang, J. and Rothrock, D.: Modeling global sea ice with a thickness and enthalpy
17 distribution model in generalized curvilinear coordinates, *Mon. Weather Rev.*, 131, 845-861,
18 2003.
- 19 Zhang, R.: Mechanisms for low-frequency variability of summer Arctic sea ice extent, *Proc.*
20 *Natl. Acad. Sci. U.S.A.*, 112, 4570-4575, doi: 10.1073/pnas.1422296112, 2015.
- 21 Zwally, H. J., Schutz, B., Abdalati, W., Abshire, J., Bentley, C., Brenner, A., Bufton, J.,
22 Dezio, J., Hancock, D., Harding, D., Herring, T., Minster, B., Quinn, K., Palm, S., Spinhirne,
23 J., and Thomas, R.: ICESat's laser measurements of polar ice, atmosphere, ocean, and land,
24 *Journal of Geodynamics*, 34, 405-445, doi: 10.1016/S0264-3707(02)00042-X, 2002.
- 25 Zygmuntowska, M., Rampal, P., Ivanova, N., and Smedsrud, L. H.: Uncertainties in Arctic
26 sea ice thickness and volume: new estimates and implications for trends, *The Cryosphere*, 8,
27 705-720, doi: 10.5194/tc-8-705-2014, 2014.

1 Table 1. List of models used: the CMIP5 subset and observations.

Institution	Model name	Ensemble members*
Commonwealth Scientific and Industrial Research Organisation (CSIRO)	CSIRO Mark version 3.6.0: CSIRO-Mk3.6.0 (Rotstayn et al., 2012)	10
Met Office Hadley Centre	Hadley Centre Global Environment Model version 2-Earth System: HadGEM2-ES (The HadGEM2 Development Team et al., 2011)	4
National Center for Atmospheric Research	Community Climate System Model, version 4: CCSM4 (Gent et al., 2011)	6
National Center for Atmospheric Research	Community Earth System Model, Community Atmosphere Model, version 5: CESM1-CAM5 (Meehl et al., 2013)	3
Model for Interdisciplinary Research on Climate (MIROC)	MIROC version 5: MIROC5 (Watanabe et al., 2010)	3
Max Plank Institute for Meteorology (MPI)	MPI Earth System Model, low resolution: MPI-ESM-LR (Jungclaus et al., 2006)	3
Applied Physics Laboratory (University of Washington)	Pan-Arctic Ice Ocean Modelling and Assimilation System: PIOMAS** (Zhang and Rothrock, 2003)	1

2 *multi-model statistics are calculated (Sect. 4.3 onwards) using the first 3 ensemble members.

3 **used as observations.

4

1 Table 2. Notation key

Notation	Description
M	Model
O_h	Observations
x_h	x over the historical period (1979 – 2014)
\bar{x}	Time mean of x over historical period
$\langle x \rangle$	Ensemble mean of x
\tilde{x}	Running time mean (11 years) of x
\hat{x}	Temporally detrended x over the historical period
σ	Standard deviation

2

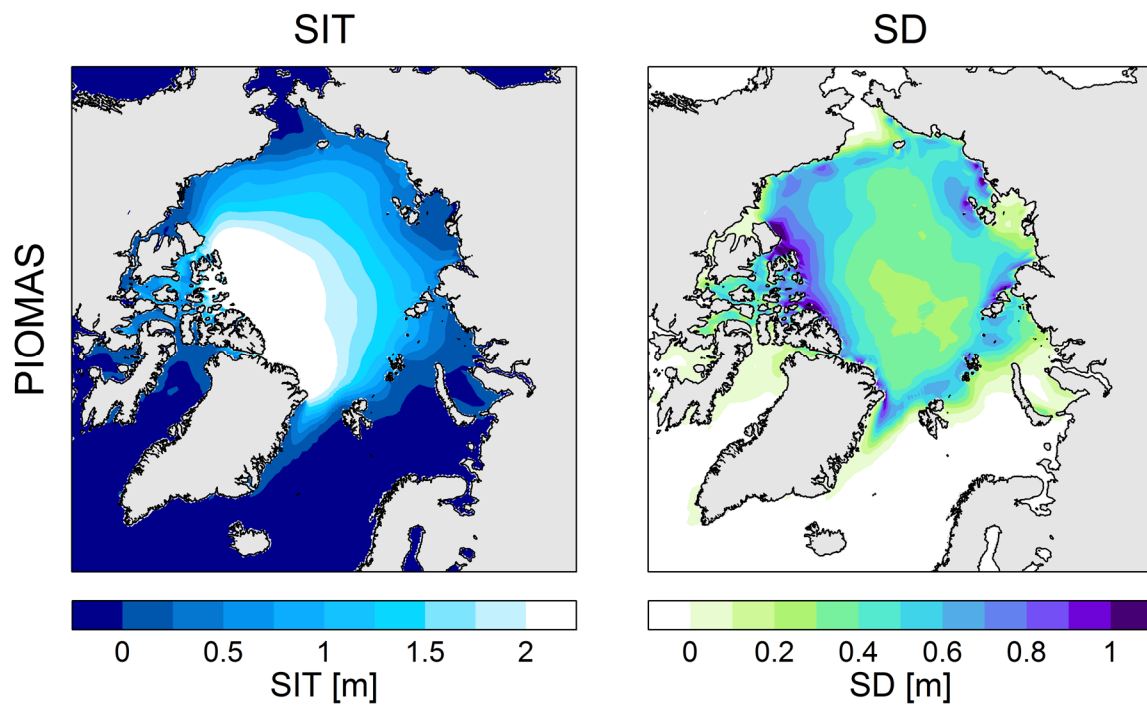


Figure 1. September 1979 – 2014 mean SIT and standard deviation (SD) from the PIOMAS reanalysis. SD is calculated after removing the linear trend.

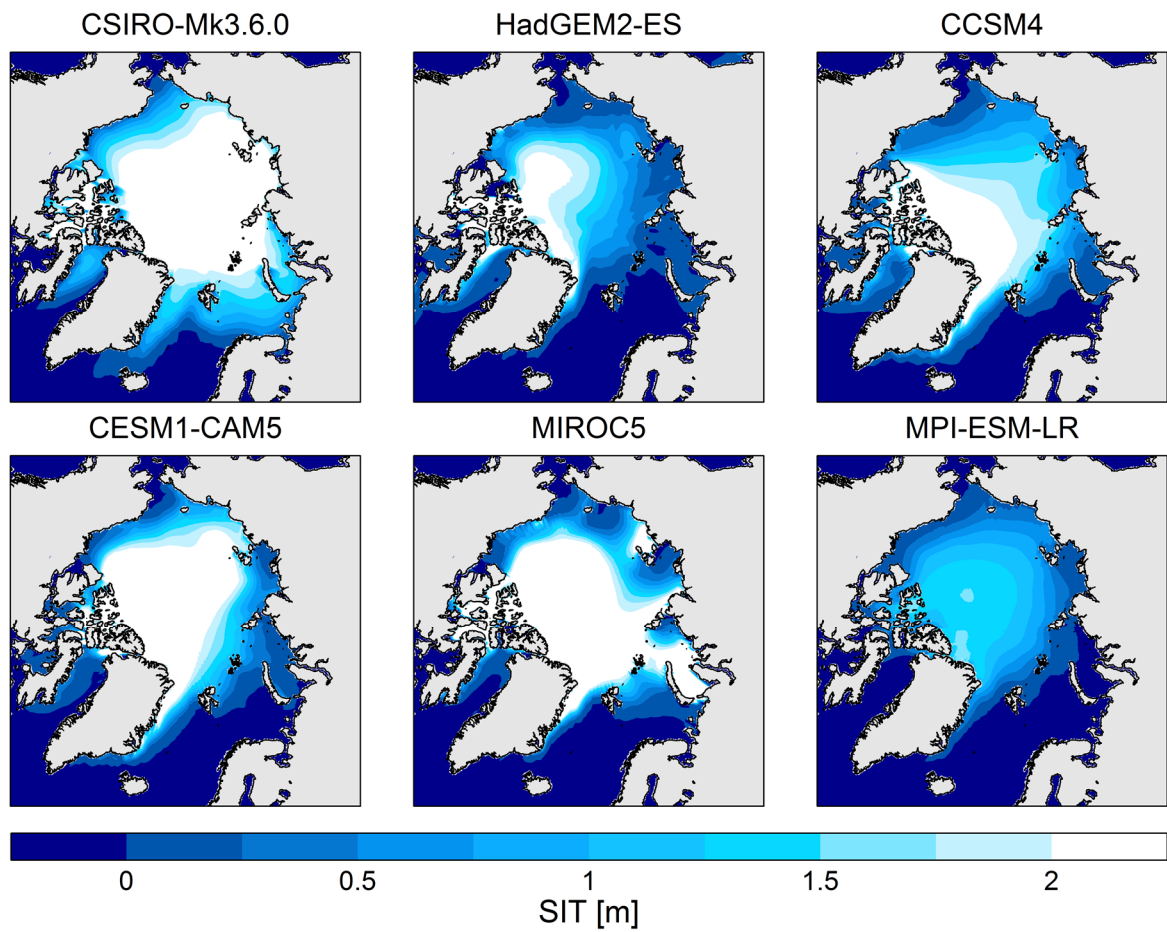


Figure 2. Mean September SIT for each of the six GCMs considered, averaged over the period 1979 – 2014.

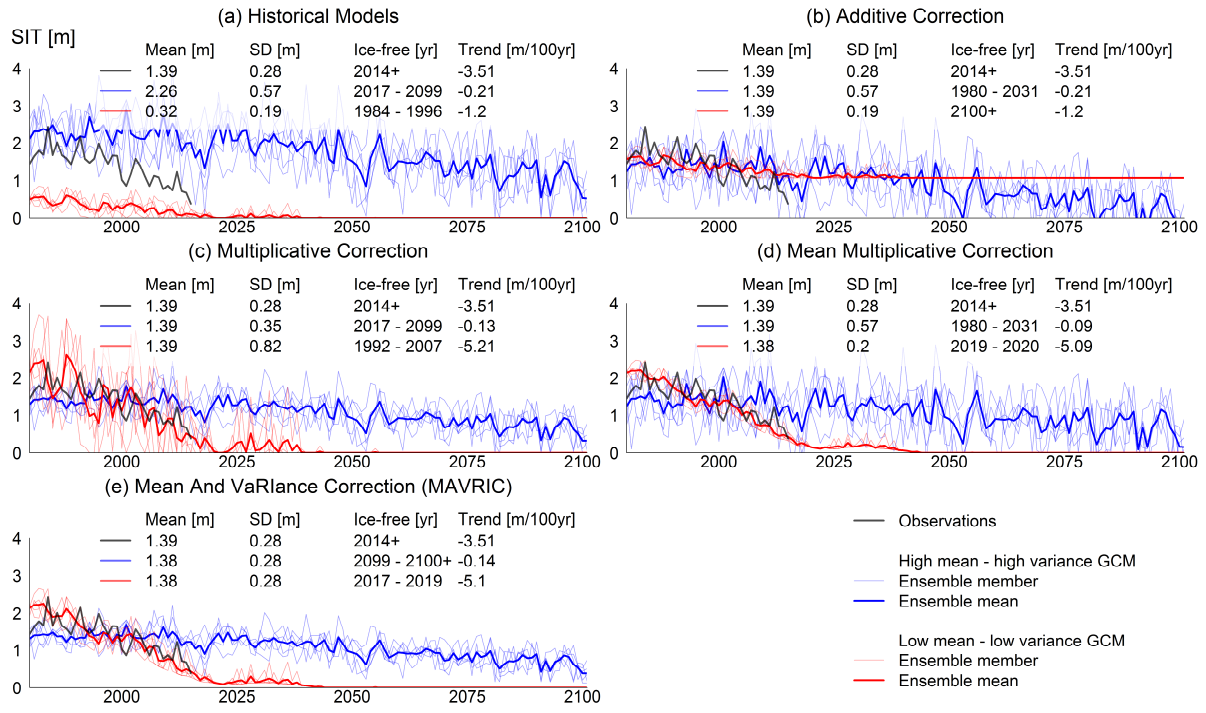


Figure 3. Performance of different SIT BCs for one particular month at a hypothetical grid point in a toy model. Mean, SD (detrended) and trend legend statistics are calculated over the observation period (1979 - 2014). ‘Ice-free’ is defined as the first occurrence of any ensemble member below 0.15 m. Shown is the ice-free ensemble range, i.e. the year of the first ensemble member to be ice-free to the last ensemble member to be ice-free. The black line represents ‘observations’, the blue and red lines represent high and low ice models respectively. The thin coloured lines represent ensemble members, and the thick lines are the ensemble mean.

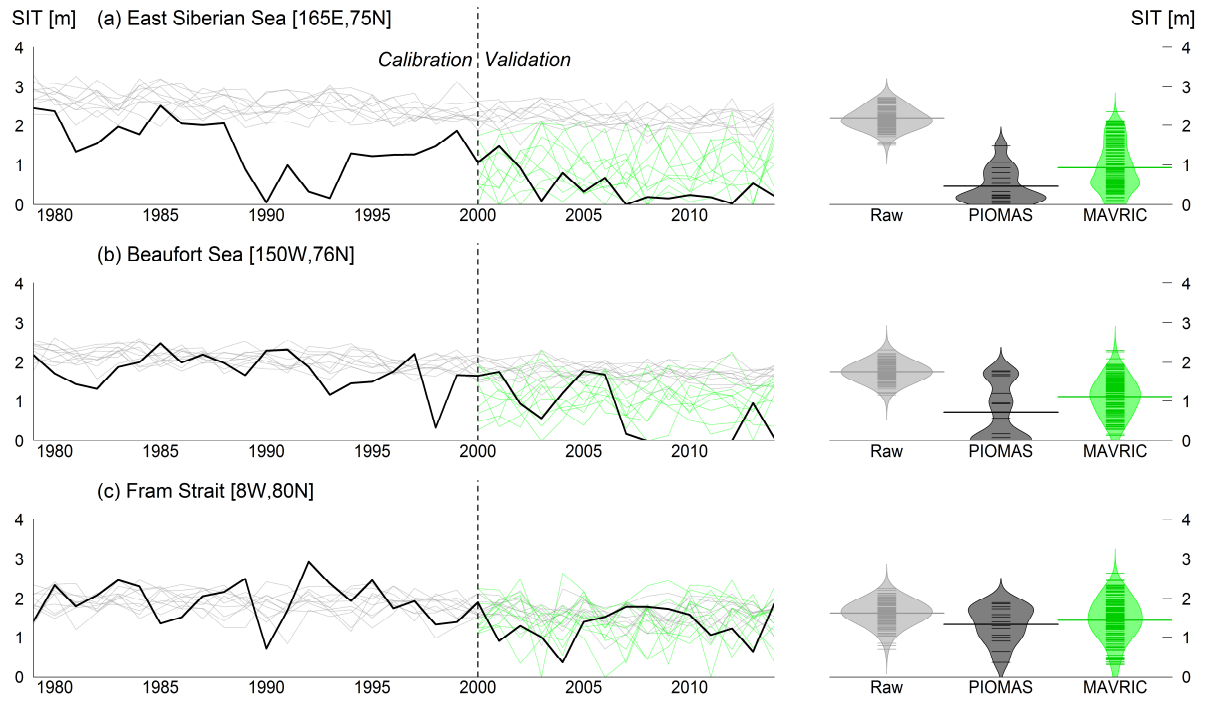


Figure 4. September SIT at three grid point locations in the Arctic, from PIOMAS (black) and CSIRO-Mk3.6.0 historical (1979 – 2005) and RCP8.5 (2006 – 2014) raw output (grey) and post MAVRIC (green). The raw CSIRO ensembles (grey) are bias corrected via the MAVRIC using the PIOMAS observations (black) over the calibration window, producing the MAVRIC ensembles (green) for the validation window. Beanplots (right) show the distribution of the SIT for the validation period. Small horizontal lines show every SIT value, the frequency of which is illustrated by the width of the shaded region. Thick horizontal line is the mean.

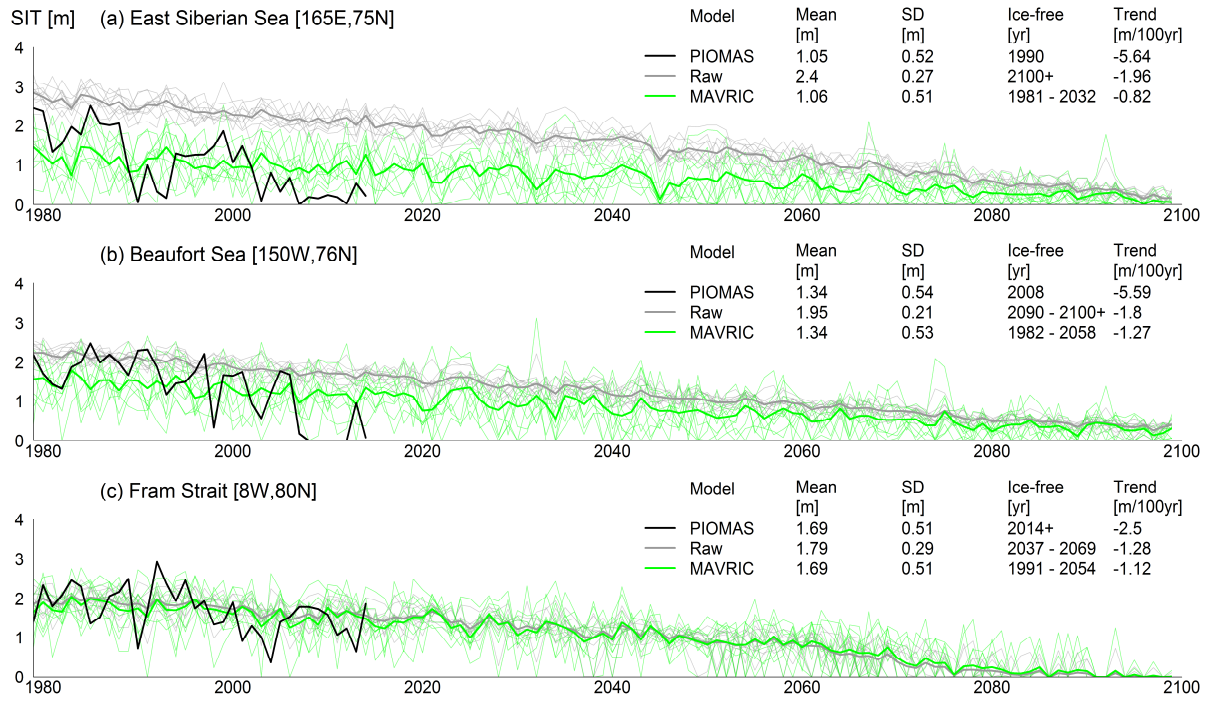


Figure 5. September SIT at three grid point locations in the Arctic, from PIOMAS (black) and CSIRO-Mk3.6.0 historical (1979 – 2005) and RCP8.5 (2006 – 2100) raw output (grey) and post MAVRIC (green). Thin lines are individual ensemble members, thick lines are the ensemble means. Mean, SD and trend legend statistics calculated over the period of observations (1979 – 2014). The SD is the detrended mean ensemble SD. Ice-free is the range of the first occurrence of the first and last ensemble member below 0.15 m.

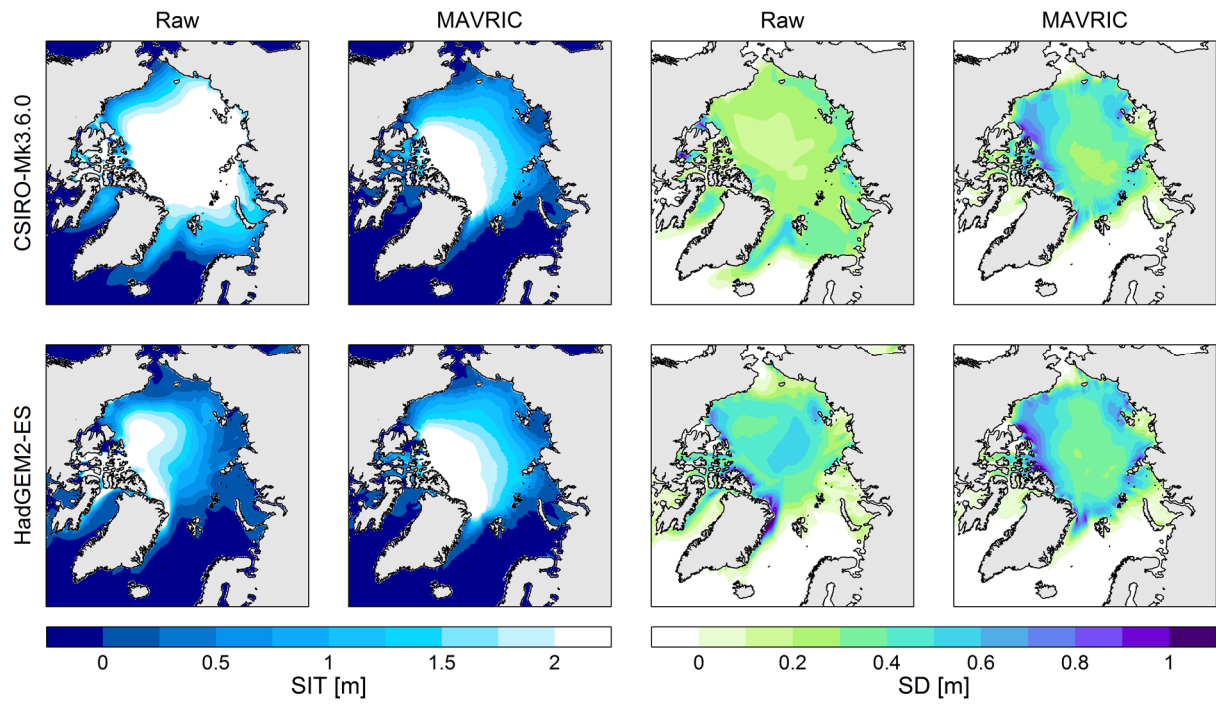


Figure 6. CSIRO-Mk3.6.0 and HadGEM2-ES, September 1979 – 2014 ensemble mean SIT and SD (detrended). The raw columns are the model solutions as found in the CMIP5 archive. The corrected columns show the distribution after the MAVRIC has been applied. PIOMAS SIT fields shown in Fig 1.

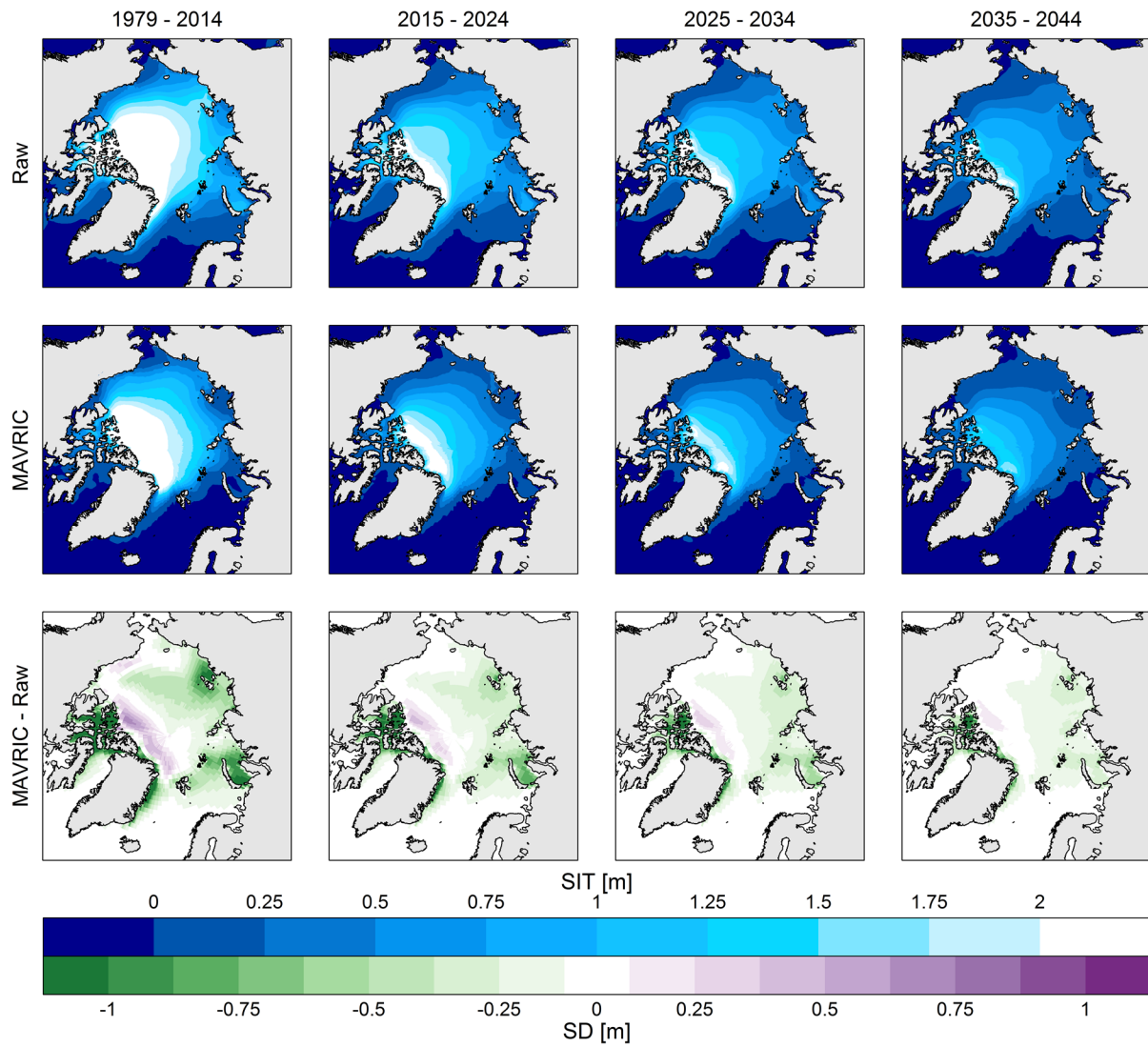


Figure 7. September multi-model ensemble mean (three members from each model) mean SIT from the CMIP5 subset, using the raw data (top row) and after MAVRIC (middle row). The bottom row shows (MAVRIC – Raw) and hence green areas are where MAVRIC has reduced SIT and purple areas are where MAVRIC has increased SIT.

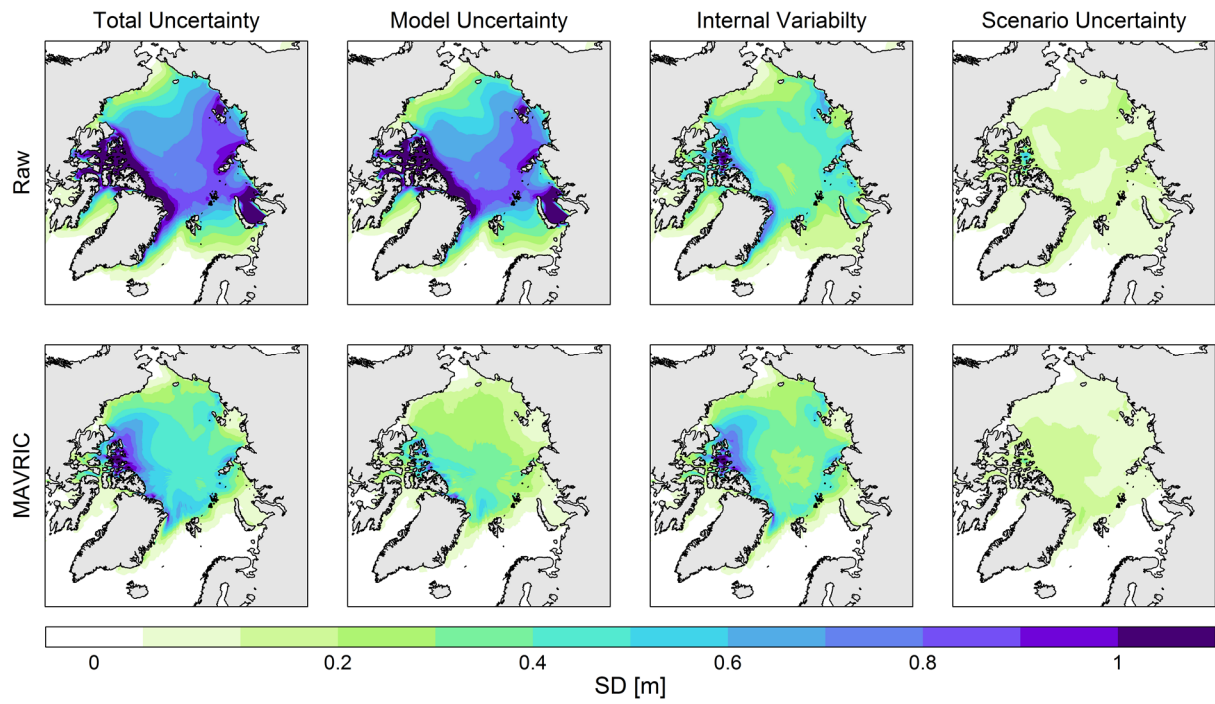


Figure 8. September 2015-2024 sources of SIT uncertainty from the CMIP5 subset (SD of the detrended SIT). The multi-model ensemble mean (three members from each) is shown when comparing raw (top row) and after MAVRIC (bottom row).

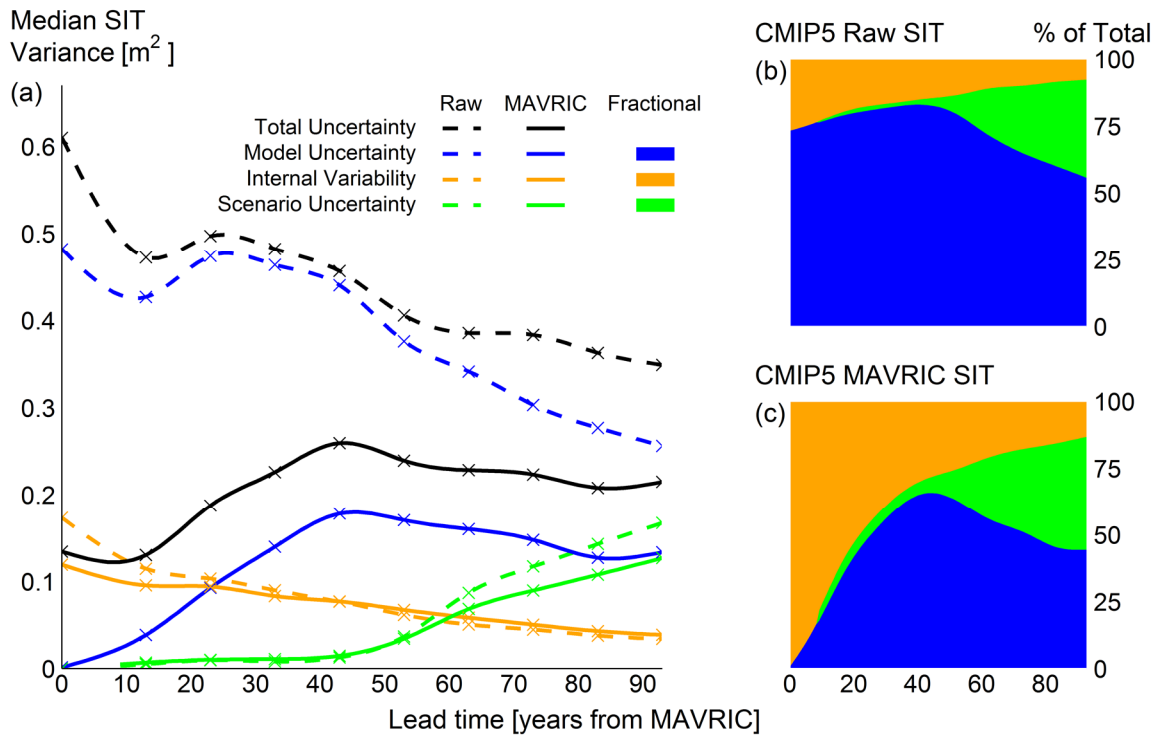


Figure 9. The evolution of the sources of September SIT uncertainty in the CMIP5 sub-set with lead time. Year zero is the MAVRIC window mid-point (1997) and the emission scenarios (RCPs) start in 2006. Panel a shows the change in magnitude of the different sources of uncertainty. The uncertainty shown is the median SIT variance and hence the lines scale additively. The dashed lines are for the raw model output and solid lines are for post MAVRIC. Contributions of model uncertainty, internal variability and scenario uncertainty as a fraction of total uncertainty are shown for the raw output (b) and post MAVRIC (c).

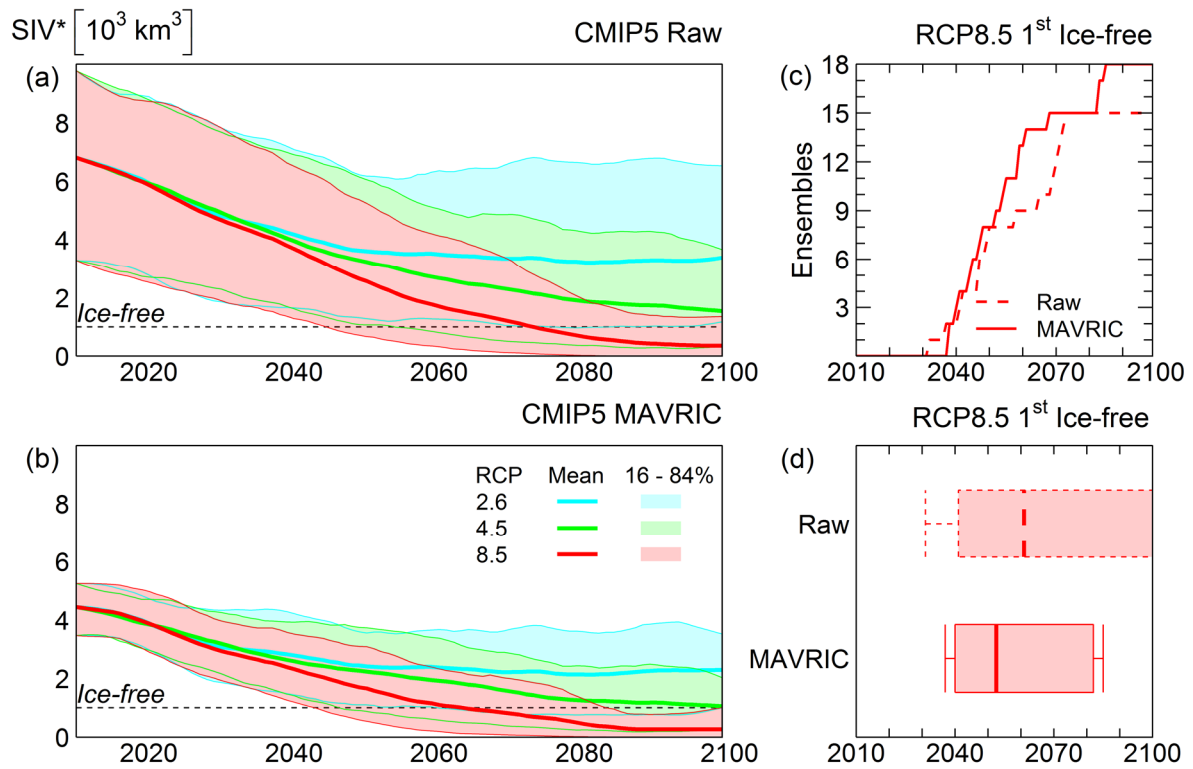
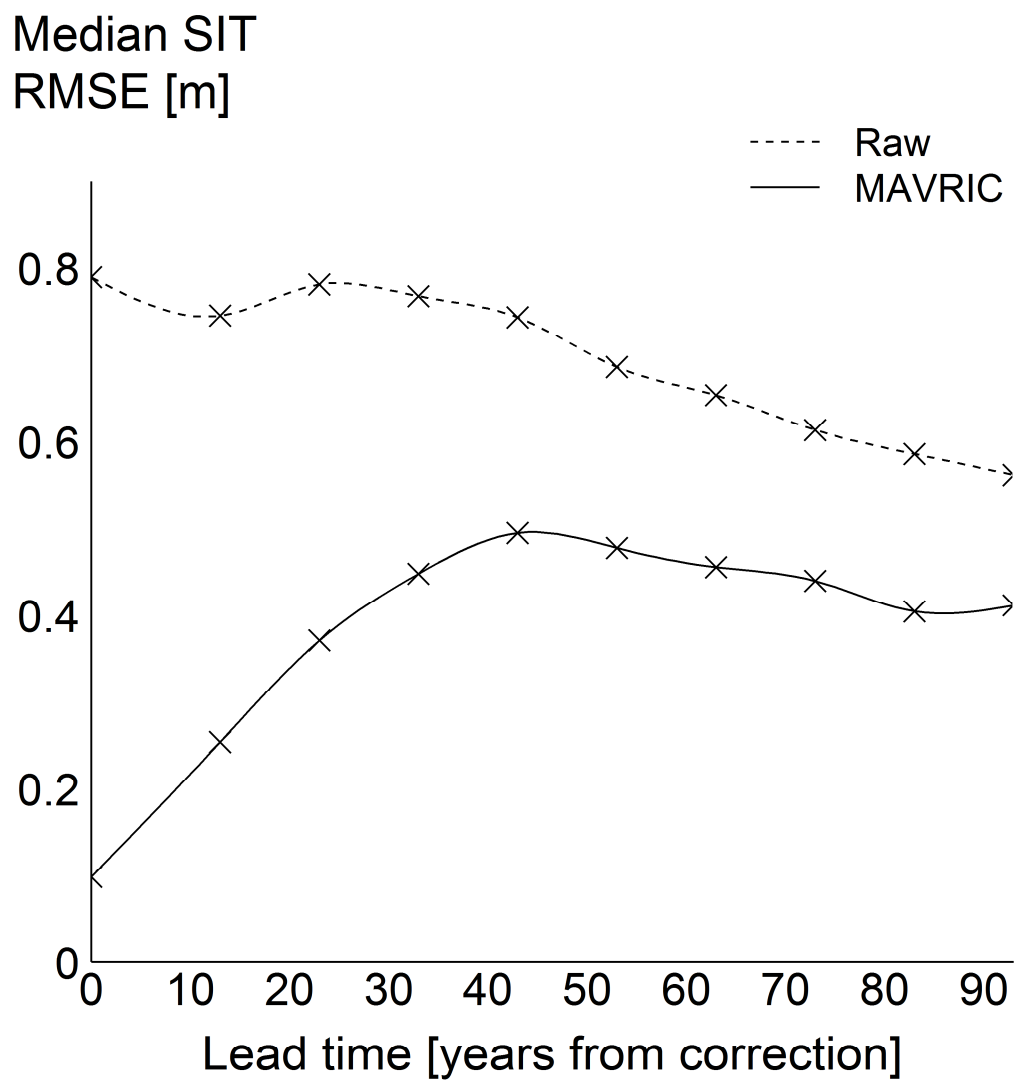
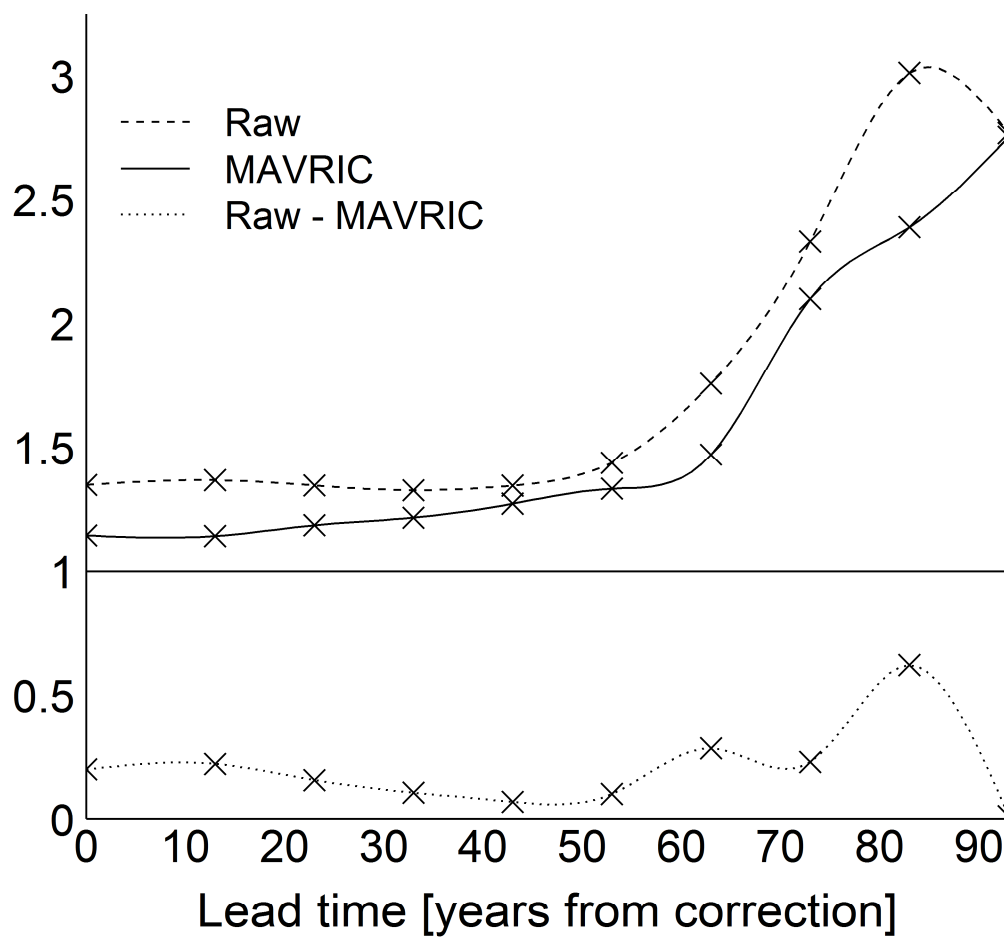


Figure 10. CMIP5 subset sea ice volume (SIV*) projections and first ice-free conditions. Panels a and b show the projected SIV* from all six models (18 ensemble members total) in both the raw and corrected GCMs (11 year running mean), and shaded regions are the 16th – 84th percentiles. Panel c shows the number of ensemble members having passed the ice-free threshold. Panel d shows the statistics of c, with the whiskers representing the range (1st and 18th ensemble member ice-free), the box capturing the 16th – 84th percentiles, and the bold line showing the median (9th ensemble member). Ice-free is defined as the first year the pan-Arctic SIV* dips below $1 \times 10^3 \text{ km}^3$ for a particular ensemble member. *Volume (SIV*) is calculated using a constant 50 % SIC throughout.



1
2 Figure C1. Multi-model ensemble out of sample September median SIT RMSE
3

Median SIT Dispersion



1
2 Figure C2. Multi-model ensemble out of sample September median SIT dispersion