

Minerva Access is the Institutional Repository of The University of Melbourne

Author/s:

Nicholls, Z;Meinshausen, M;Lewis, J;Corradi, MR;Dorheim, K;Gasser, T;Gieseke, R;Hope, AP;Leach, NJ;McBride, LA;Quilcaille, Y;Rogelj, J;Salawitch, RJ;Samset, BH;Sandstad, M;Shiklomanov, A;Skeie, RB;Smith, CJ;Smith, SJ;Su, X;Tsutsui, J;Vega-Westhoff, B;Woodard, DL

Title:

Reduced Complexity Model Intercomparison Project Phase 2: Synthesizing Earth System Knowledge for Probabilistic Climate Projections

Date:

2021-06-01

Citation:

Nicholls, Z., Meinshausen, M., Lewis, J., Corradi, M. R., Dorheim, K., Gasser, T., Gieseke, R., Hope, A. P., Leach, N. J., McBride, L. A., Quilcaille, Y., Rogelj, J., Salawitch, R. J., Samset, B. H., Sandstad, M., Shiklomanov, A., Skeie, R. B., Smith, C. J., Smith, S. J., ... Woodard, D. L. (2021). Reduced Complexity Model Intercomparison Project Phase 2: Synthesizing Earth System Knowledge for Probabilistic Climate Projections. *Earth S Future*, 9 (6), <https://doi.org/10.1029/2020EF001900>.

Persistent Link:

<https://hdl.handle.net/11343/280594>

License:

[CC BY](#)

# Earth's Future



## RESEARCH ARTICLE

10.1029/2020EF001900

### Key Points:

- Probabilistic global-mean temperature projections often use reduced complexity climate models (RCMs) because of their low computational cost
- We evaluate how well RCMs' probabilistic setups can synthesize knowledge from multiple research domains for policy relevant projections
- No RCM is able to capture all forcing, warming, heat uptake and carbon cycle metrics we evaluate, however some come close across a range

### Supporting Information:

Supporting Information may be found in the online version of this article.

### Correspondence to:

Z. Nicholls,  
[zebedee.nicholls@climate-energy-college.org](mailto:zebedee.nicholls@climate-energy-college.org)

### Citation:

Nicholls, Z., Meinshausen, M., Lewis, J., Corradi, M. R., Dorheim, K., Gasser, T., et al. (2021). Reduced complexity Model Intercomparison Project Phase 2: Synthesizing Earth system knowledge for probabilistic climate projections. *Earth's Future*, 9, e2020EF001900. <https://doi.org/10.1029/2020EF001900>









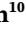






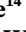
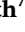





Received 12 NOV 2020

Accepted 25 APR 2021

### Author Contributions:

**Conceptualization:** Z. Nicholls, M. Meinshausen, M. Rojas Corradi, K. Dorheim, T. Gasser, R. Gieseke, A. P. Hope, N. J. Leach, L. A. McBride, Y. Quilcaille, J. Rogelj, R. J. Salawitch, B. H. Samset, M. Sandstad, A. Shiklomanov, R. B. Skeie, C. J. Smith,

## Reduced Complexity Model Intercomparison Project Phase 2: Synthesizing Earth System Knowledge for Probabilistic Climate Projections

Z. Nicholls<sup>1,2</sup> , M. Meinshausen<sup>1,2,3</sup> , J. Lewis<sup>1</sup> , M. Rojas Corradi<sup>4,5</sup> , K. Dorheim<sup>6</sup> , T. Gasser<sup>7</sup> , R. Gieseke<sup>8</sup> , A. P. Hope<sup>9</sup> , N. J. Leach<sup>10</sup> , L. A. McBride<sup>11</sup> , Y. Quilcaille<sup>7</sup> , J. Rogelj<sup>7,12</sup> , R. J. Salawitch<sup>9,11,13</sup> , B. H. Samset<sup>14</sup> , M. Sandstad<sup>14</sup>, A. Shiklomanov<sup>15</sup> , R. B. Skeie<sup>14</sup> , C. J. Smith<sup>7,16</sup> , S. J. Smith<sup>17</sup> , X. Su<sup>18</sup> , J. Tsutsui<sup>19</sup> , B. Vega-Westhoff<sup>20</sup> , and D. L. Woodard<sup>5</sup> 

<sup>1</sup>Australian-German Climate & Energy College, University of Melbourne, Parkville, VIC, Australia, <sup>2</sup>School of Earth Sciences, University of Melbourne, Parkville, VIC, Australia, <sup>3</sup>Potsdam Institute for Climate Impact Research, Member of the Leibniz Association, Potsdam, Germany, <sup>4</sup>Department of Geophysics, University of Chile, Santiago, Chile, <sup>5</sup>Center for Climate and Resilience Research, CR2, Santiago, Chile, <sup>6</sup>Pacific Northwest National Laboratory, Richland, WA, USA, <sup>7</sup>International Institute for Applied Systems Analysis, Laxenburg, Austria, <sup>8</sup>Independent Researcher, Potsdam, Germany, <sup>9</sup>Department of Atmospheric and Oceanic Science, University of Maryland-College Park, College Park, USA, <sup>10</sup>Department of Physics, Atmospheric, Oceanic, and Planetary Physics, University of Oxford, Oxford, UK, <sup>11</sup>Department of Chemistry and Biochemistry, University of Maryland-College Park, College Park, MD, USA, <sup>12</sup>Grantham Institute, Imperial College London, London, UK, <sup>13</sup>Earth System Science Interdisciplinary Center, University of Maryland-College Park, College Park, MD, USA, <sup>14</sup>CICERO Center for International Climate Research, Oslo, Norway, <sup>15</sup>NASA Goddard Space Flight Center, Greenbelt, MD, USA, <sup>16</sup>Priestley International Centre for Climate, University of Leeds, Leeds, UK, <sup>17</sup>Joint Global Change Research Institute, Pacific Northwest National Laboratory, College Park, MD, USA, <sup>18</sup>Research Institute for Global Change/Research Center for Environmental Modeling and Application/Earth System Model Development and Application Group, Japan Agency for Marine-Earth Science and Technology, Yokohama, Japan, <sup>19</sup>Environmental Science Research Laboratory, Central Research Institute of Electric Power Industry, Abiko, Japan, <sup>20</sup>Department of Atmospheric Sciences, University of Illinois at Urbana-Champaign, Urbana, IL, USA

**Abstract** Over the last decades, climate science has evolved rapidly across multiple expert domains. Our best tools to capture state-of-the-art knowledge in an internally self-consistent modeling framework are the increasingly complex fully coupled Earth System Models (ESMs). However, computational limitations and the structural rigidity of ESMs mean that the full range of uncertainties across multiple domains are difficult to capture with ESMs alone. The tools of choice are instead more computationally efficient reduced complexity models (RCMs), which are structurally flexible and can span the response dynamics across a range of domain-specific models and ESM experiments. Here we present Phase 2 of the Reduced Complexity Model Intercomparison Project (RCMIP Phase 2), the first comprehensive intercomparison of RCMs that are probabilistically calibrated with key benchmark ranges from specialized research communities. Unsurprisingly, but crucially, we find that models which have been constrained to reflect the key benchmarks better reflect the key benchmarks. Under the low-emissions SSP1-1.9 scenario, across the RCMs, median peak warming projections range from 1.3 to 1.7°C (relative to 1850–1900, using an observationally based historical warming estimate of 0.8°C between 1850–1900 and 1995–2014). Further developing methodologies to constrain these projection uncertainties seems paramount given the international community's goal to contain warming to below 1.5°C above preindustrial in the long-term. Our findings suggest that users of RCMs should carefully evaluate their RCM, specifically its skill against key benchmarks and consider the need to include projections benchmarks either from ESM results or other assessments to reduce divergence in future projections.

**Plain Language Summary** Our best tools to capture state-of-the-art knowledge are complex, fully coupled Earth System Models (ESMs). However, ESMs are expensive to run and no single ESM can easily produce responses which represent the full range of uncertainties. Instead, for some applications, computationally efficient reduced complexity climate models (RCMs) are used in a probabilistic setup. An example of these applications is estimating the likelihood that an emissions scenario will stay below a certain global-mean temperature change. Here we present a study (referred to as the Reduced Complexity

© 2021. The Authors. Earth's Future published by Wiley Periodicals LLC on behalf of American Geophysical Union. This is an open access article under the terms of the [Creative Commons Attribution License](https://creativecommons.org/licenses/by/4.0/), which permits use, distribution and reproduction in any medium, provided the original work is properly cited.

S. J. Smith, X. Su, J. Tsutsui, B. Vega-Westhoff, D. L. Woodard  
**Data curation:** J. Lewis, K. Dorheim, T. Gasser, R. Gieseke, A. P. Hope, N. J. Leach, L. A. McBride, Y. Quilcaille, R. J. Salawitch, B. H. Samset, M. Sandstad, A. Shiklomanov, R. B. Skeie, C. J. Smith, S. J. Smith, X. Su, J. Tsutsui, B. Vega-Westhoff, D. L. Woodard  
**Methodology:** Z. Nicholls, M. Meinshausen, K. Dorheim, T. Gasser, R. Gieseke, A. P. Hope, N. J. Leach, L. A. McBride, Y. Quilcaille, J. Rogelj, R. J. Salawitch, B. H. Samset, M. Sandstad, A. Shiklomanov, R. B. Skeie, C. J. Smith, S. J. Smith, X. Su, J. Tsutsui, B. Vega-Westhoff, D. L. Woodard  
**Software:** Z. Nicholls, J. Lewis, R. Gieseke  
**Supervision:** M. Meinshausen  
**Validation:** J. Lewis  
**Visualization:** Z. Nicholls  
**Writing – original draft:** Z. Nicholls  
**Writing – review & editing:** Z. Nicholls, M. Meinshausen, J. Lewis, M. Rojas Corradi, K. Dorheim, T. Gasser, R. Gieseke, A. P. Hope, N. J. Leach, L. A. McBride, Y. Quilcaille, J. Rogelj, R. J. Salawitch, B. H. Samset, M. Sandstad, A. Shiklomanov, R. B. Skeie, C. J. Smith, S. J. Smith, X. Su, J. Tsutsui, B. Vega-Westhoff, D. L. Woodard

Model Intercomparison Project (RCMIP) Phase 2) which investigates the extent to which different RCMs can be probabilistically calibrated to reproduce knowledge from specialized research communities. We find that the agreement between each RCM and the benchmarks varies, although the best performing models show good agreement across the majority of benchmarks. Under a very-low-emissions scenario median peak warming projections range from 1.3 to 1.7°C (relative to 1850–1900, assuming historical warming of 0.8°C between 1850–1900 and 1995–2014). Investigating new ways to reduce these projection uncertainties seems paramount given the international community's goal to limit warming to below 1.5°C above preindustrial in the long-term.

## 1. Introduction

Coupled Earth System Models (ESMs) have evolved for decades as primary climate research tools (Kawamiya et al., 2020). They represent the state of the art of complex Earth system modeling. Nonetheless, they are not the tool of choice to assess the full breadth of scenario and Earth system response uncertainty that has been identified in the scientific literature. It is infeasible to assess the climate implications of hundreds to thousands of emissions scenarios with the world's most comprehensive ESMs, such as those participating in the Sixth Phase of the Couple Model Intercomparison Project (CMIP6) (Eyring et al., 2016), because of ESMs' computational cost, the complexity in setting up input data and the sheer volume of output data generated. Yet, large scenario assessments are vital for understanding the consequences of various policy choices and their residual climate hazards.

Similarly, while some ESMs perform large, perturbed physics experiments (e.g., Murphy et al., 2014) that aim to explore a range of potential Earth system long-term annual-average responses, the ability to capture full uncertainty ranges is limited. The ability to capture full uncertainty ranges is limited because these ESMs are relatively rigid in their structure—lacking the ability to completely explore uncertainties in vital components like the carbon cycle or effective radiative forcings.

An answer to both of these challenges, i.e., (a) limited computational resources and (b) structural scope and flexibility to represent long-term uncertainties in key metrics like global-mean surface air temperatures, are Reduced Complexity Models (RCMs), often also referred to as simple climate models (SCMs). RCMs can play the vital role of extending the knowledge and uncertainties from multiple domains, particularly a multitude of ESM experiments, to probabilistic long-term climate projections of key variables over a wide range of scenarios. Earth System Models of Intermediate Complexity (EMICs) may initially appear to be another option. However, due to the process-based representations used by EMICs, their computational complexity and data requirements are still orders of magnitude greater than RCMs. As a result, even EMICs are not a feasible choice for the large-scale, probabilistic assessment discussed here.

Typically, RCMs achieve computational efficiency and structural flexibility by limiting their spatial and temporal domains to global-mean, annual-mean quantities, i.e., the domains of relevance to long-term, global climate change. In general, RCMs do not include representations of interannual variability, although the EMGC model (Table 1) is a clear exception to this rule. Rather than aiming to represent the physics of the climate system at the process level and high-resolution, RCMs use parameterizations of the system which capture its large-scale behavior at a greatly reduced computational cost. This allows them to perform 350-year long simulations in a fraction of a second on a single CPU, multiple orders of magnitude faster than our most comprehensive ESMs which would take weeks to months on the world's most advanced supercomputers.

A key example of large-scale emissions scenario assessment, and the one we focus on in this paper, is the climate assessment of socioeconomic scenarios by Working Group 3 (WG3) of the Intergovernmental Panel on Climate Change (IPCC). Hundreds of emission scenarios were assessed in the IPCC's Fifth Assessment Report (AR5, see Clarke et al. (2014)) as well as its more recent Special Report on Global Warming of 1.5°C (SR1.5, see Rogelj, Shindell, et al. (2018) and Huppmann et al. (2018)) (Scenario data is available at <https://secure.iiasa.ac.at/web-apps/ene/AR5DB> and <https://data.ene.iiasa.ac.at/iamc-1.5c-explorer/> for AR5 and SR1.5 respectively, both databases are hosted by the IIASA Energy Program). For the IPCC's forthcoming Sixth Assessment (AR6), it is anticipated that the number of scenarios will be in the several hundreds to a thousand (an initial snapshot of scenarios based on the SSPs is available at <https://tntcat.iiasa.ac.at/SspDb>).

**Table 1**  
Overview of the Models and Constraining Approaches Used in this Paper

Model	Constraining technique	Key references
CICERO-SCM	591 members subsampled from a posterior of 30,040 members to form a set that match the proxy assessment ocean heat content distribution while excluding parameter sets with unrealistic aerosol ERF or unrealistic surface air temperature change from 1850–1900 to 1985–2014	Schlesinger et al. (1992), Joos et al. (1996), Etminan et al. (2016), Skeie et al. (2017, 2018), Z. R. J. Nicholls et al. (2020), and Skeie et al. (2021)
EMGC	160,000 sample members, retaining the 1,000 that minimize reduced-chi-squared between modeled and observed GMST and OHC from 1850 to 1999	Canty et al. (2013), Hope et al. (2017, 2020), and McBride et al. (2021)
FaIRv1.6.1	3,000 sample members retaining the 501 that minimize RMSE between modeled and observed 1850–2014 GMST	Millar et al. (2017) and Smith, Forster, et al. (2018)
FaIRv2.0.0-alpha	1,000,000 member raw ensemble, constrained with likelihood of 2010–2019 level and rate of attributable warming, calculated using the Global Warming Index methodology (Haustein et al., 2017). 5000 members randomly drawn from the constrained ensemble for use here.	Millar et al. (2017), Haustein et al. (2017), Smith, Forster, et al. (2018), and Leach et al. (2020)
Hectorv2.5.0	10,000 member ensemble sampled from Markov chain Monte Carlo chains constrained with global surface temperature and ocean heat content	Vega-Westhoff et al. (2019)
MAGICCv7.5.1	7,000,000 member Monte Carlo Markov Chain, 600 member subsample selected to match proxy assessed ranges	Meinshausen et al. (2009, 2011, 2020)
MCE v1.2	600 members sampled with a Metropolis-Hastings algorithm through Bayesian updating to reflect an ensemble of complex climate models constrained with the proxy assessed ranges	Tsutsui (2017, 2020) (see also Joos et al. (1996) and Hooss et al. (2001))
OSCARv3.1	10,000 Monte Carlo members, weighted using their agreement with a set of assessed ranges (supplementary Text S1)	Gasser et al. (2017, (2018, 2020)
SCM4OPT v2.1	For each emission scenario, 2,000 sample members are used to reflect uncertainties resulting from carbon cycle, aerosol forcings and temperature change, while constrained by the historical mean surface temperature of HadCRUT.4.6.0.0 (Morice et al., 2012)	Su et al. (2017, 2018, 2020)

Note. Detailed descriptions of each model are available in supplementary Text S1.

Running WG3-type scenarios requires at least some representation of greenhouse gas cycles, atmospheric chemistry and dynamic vegetation modules. While some of the world's most comprehensive ESMs have the required components, they could not be used to sample scenario and parametric uncertainty for reasons of computational cost. The most comprehensive RCMs include parameterized representations of the required components (including feedbacks of climate on permafrost and wetland methane emissions), enabling the exploration of interacting uncertainties from multiple parts of the climate system in an internally consistent setup.

While RCMs do not include the detail of ESMs across the emissions-climate change cause-effect chain, they do tend to include uncertainty representations for more steps in the chain (i.e., RCMs tradeoff depth for breadth compared to ESMs). For example, many RCMs include the relationship between methane emissions and concentrations (including temperature and other feedbacks) whereas few ESMs do in their long-term experiments. On the other hand, few RCMs directly use land-cover information within their carbon cycles, and none consider it in the detailed way which ESMs do. In addition, there are clearly applications where RCMs are not a feasible tool. For example, near-term attribution studies, such as the World Weather Attribution project (Uhe et al., 2016). For this latter application, large-ensemble ESM runs are vital—as only they can reflect natural variability and weather patterns. Overall, there is no question that ESMs are by far the most important research tool to project future climate change. RCMs complement the ESM efforts. Within this paper, we focus on a very specific niche of this complementing role, i.e., the degree to which RCMs can synthesize multiple lines of evidence across the emissions-climate change cause-effect chain.

Typically, RCMs attempt to perform this synthesis using probabilistic parameter ensembles (see also Section 3), which are distinct from the emulator mode in which RCMs can also be run (see Nicholls et al., 2020 for a discussion of emulation with RCMs). These probabilistic parameter ensembles are derived based on knowledge of specific Earth system quantities drawn from multiple, often independent, research

communities, e.g., historical global-mean temperature increase, effective radiative forcing due to different anthropogenic emissions, ocean heat uptake, or cumulative land and ocean carbon uptake. The resulting distributions can then be used in a variety of applications, e.g., to assess the likelihood that different warming levels are reached under a specific emissions scenario (e.g., 50% and 66%) based on the combined available evidence. As a result of their probabilistic nature, the ensembles resulting from RCMs are conceptually different from an ensemble of multiple model outputs that has not been constructed with any relative probabilities in mind (such as those from CMIP6) taken without constraining or any other sort of postprocessing.

Within the IPCC, RCMs' synthesizing niche facilitates the transfer of knowledge from Working Group I (WG1), which assesses the physical science of the climate system, to WG3, which assesses the socioeconomics of climate change mitigation. The goal of this knowledge transfer is consistency between WG3's scenario classification and the physical science assessment of WG1—a key precondition to have confidence that WG3's conclusions about the socioeconomic transformation required to mitigate anthropogenic climate change to specific levels are based on our latest scientific understanding. Here, we describe RCMs as “integrators of knowledge” because they integrate (a relevant subsection of) the assessment from WG1, providing WG3 with a tool that can be used for assessing the climate implications, particularly global-mean temperature changes, of a wide range of emissions scenarios.

Due to their role in the IPCC assessment (and for analyzing mitigation options in line with temperature targets more generally), understanding the degree to which RCMs can reflect a range of independent radiative forcing, warming, heat uptake and concentration assessments simultaneously is of vital importance. Given that these assessments are independent, a single, internally consistent, model may not be able to capture them all. If RCMs are inherently biased in some way or they are unable to simultaneously capture the independent assessments, this will affect the WG3 climate assessment and interpretation of the RCMs' outputs should be adjusted accordingly.

This study's scope, in terms of number of climate dimensions considered and number of climate models evaluated, is unique. While there have been studies with single models which choose parameter sets that match various assessments of ECS and TCR (e.g., Meinshausen, Meinshausen, et al., 2009; Rogelj, Meinshausen, & Knutti, 2012) and Smith, Forster, et al. (2018) compared two models' probabilistic outputs, no previous study into RCM probabilistic distributions is of the same breadth.

Here, in the second phase of RCMIP, we evaluate the degree to which multiple RCMs are able to synthesize Earth system knowledge within a probabilistic distribution. We then examine the implications of differences in these probabilistic distributions for climate projections. We extend previous probabilistic evaluation work and build on the progress made in the first phase of RCMIP (Z. R. J. Nicholls et al., 2020) and other RCM intercomparison studies (Harmsen et al., 2015; Schwarber et al., 2019; van Vuuren et al., 2011). We widen the first phase's scope both in terms of number of climate dimensions considered and the number of models evaluated. To our knowledge, this is the most comprehensive evaluation performed to date of the ability of RCMs to capture a broad range of climate metrics and key indicators, such as those assessed by IPCC WG1.

## 2. Participating Models

Nine models have participated in RCMIP Phase 2 (Table 1 and supplementary Text S1). Models were invited to participate via an open invitation made available at [rcmip.org](http://rcmip.org) and circulated via relevant researcher networks. All interested modeling teams were included. These models and their components range from simpler, regression-based approaches to more complex representations with detailed processes and regions. The models have been constrained in a number of different ways, using statistical techniques ranging in complexity from Monte Carlo Markov Chains to using pass/fail criteria to determine valid parameter values. As a result, the models and techniques cover (to the best of our knowledge) the full range of techniques seen in the literature and their results allow us to evaluate the implications of different choices.

**Table 2**  
*The Proxy Assessed Ranges Used in this Study*

Metric	Assessed range Unit	vll	ll	c	lu	vlu
2000–2019 GMST rel. to 1961–1990	K	0.46		0.54		0.61
Equilibrium climate sensitivity	K	2.30	2.60	3.10	3.90	4.70
Transient climate response	K	0.98	1.26	1.64	2.02	2.29
Transient climate response to emissions	K/TtC	1.03	1.40	1.77	2.14	2.51
2014 CO <sub>2</sub> effective radiative forcing	W/m <sup>2</sup>		1.69	1.80	1.91	
2014 Aerosol effective radiative forcing	W/m <sup>2</sup>		−1.37	−1.01	−0.63	
2018 Ocean heat content rel. to 1971	ZJ		303	320	337	
2011 CH <sub>4</sub> effective radiative forcing	W/m <sup>2</sup>		0.47	0.60	0.73	
2011 N <sub>2</sub> O effective radiative forcing	W/m <sup>2</sup>		0.14	0.17	0.20	
2011 F-Gases effective radiative forcing	W/m <sup>2</sup>		0.03	0.03	0.03	

*Note.* The assessed ranges are labeled as “vll” (very likely lower, i.e., Fifth percentile), “ll” (likely lower, 17th percentile), “c” (central, 50th percentile), “lu” (likely upper, 83rd percentile), and “vlu” (very likely upper, 95th percentile). Sources are described in Section 3.

### 3. Methods

In this study, the RCMs are run in a probabilistic setup, also referred to as a probabilistic distribution. As discussed in Section 1, a probabilistic setup means that each RCM is run with an ensemble of parameter configurations. Specifically, for a given experiment, each RCM is run multiple times, each time in a different configuration, i.e., with different parameter values. All of these different runs are then combined to form a probabilistic set of outputs. With these probabilistic sets, we can then calculate ranges of each output variable of interest (e.g., global-mean surface temperatures).

Modeling groups use a range of techniques to derive their parameter ensembles, i.e., to constrain their models (Table 1). In each probabilistic run, the parameter ensemble is fixed, i.e., the same set of parameter configurations will be used in each experiment. This choice ensures that the model outputs are deterministic, rather than including a random element due to, e.g., sampling parameter values from a range or probability distribution for each run. Typically, modeling groups will also use different data to derive their parameter ensemble. This can lead to differences in model projections which are simply based on choices made by the modeling groups and are not related to model structure or constraining technique at all. In this study, two models (MAGICC7 and MCE-v1-2) have used a common set of target assessed ranges, i.e., benchmarks, to derive their probabilistic distributions. For these models, we are able to rule out the choice of data as the cause of difference between these models. Accordingly, we can more clearly identify the importance of model structure and constraining technique for future projections.

In this study, our target assessment is a “proxy assessment”, which uses assessed climate system characteristics in line with IPCC AR5 as its starting point and updates key values using more recent literature (Table 2). We explicitly use the name “proxy assessment” throughout to make clear that we are not constraining to any ranges coming from the formal IPCC assessment, rather an approximation thereof. Notably, in this study, the proxy assessment does not include any future projections. While we examine future projections coming from the models, we do not explicitly compare them against future projections coming from another line of evidence because there is no obvious choice for such a line of evidence—apart from the “assessed ranges” of SSP scenarios that will be communicated in the forthcoming IPCC report (but are not available for this study). As discussed in more detail in Section 4.3, the inclusion of future projections in the proxy assessment would narrow the range of model projections but any such narrowing should be carefully considered because—depending on the types of constraints—it may lead to underestimates of uncertainty.

In order to keep the study’s scope manageable, our proxy assessment focuses on climate response parameters, with the carbon cycle examined only via the TCRE. We aim to perform a detailed analysis on carbon cycle response in the next phase of RCMIP.

We use surface air ocean blended temperatures from the HadCRUT.4.6.0.0 data set (Morice et al., 2012). HadCRUT4.6.0.0 is a widely used observational data product and is representative of other observations of changes in surface air and ocean temperatures (Simmons et al., 2017). Our key metric for evaluating RCM temperature projections is the warming between the 1961–1990 and 2000–2019 periods (using the SSP2-4.5 scenario to extend the CMIP6 historical experiment to 2019). We choose a relatively recent period to match the increase in global observations since the 1960s.

For ocean heat content, we use the recent work of von Schuckmann et al. (2020). We focus on the change in ocean heat content between 1971 and 2018, when the largest set of observations are available.

We use the recent assessment of Sherwood et al. (2020) for equilibrium climate sensitivity (ECS). ECS is defined as the equilibrium warming which occurs under a doubling of atmospheric CO<sub>2</sub> concentrations relative to preindustrial concentrations. The ECS assessment is combined with the constrained transient climate response (TCR) assessment of Tokarska et al. (2020). TCR is defined as the surface air temperature change which occurs at the time at which atmospheric CO<sub>2</sub> concentrations double in an experiment in which atmospheric CO<sub>2</sub> concentrations rise at one percent per year (a 1pctCO<sub>2</sub> experiment). Carbon cycle behavior is considered only via the transient climate response to emissions (TCRE). TCRE is defined as the ratio of surface air temperature change to cumulative CO<sub>2</sub> emissions at the time when atmospheric CO<sub>2</sub> concentrations double in a 1pctCO<sub>2</sub> experiment. We use the TCRE assessment from Arora, Katavouta, et al. (2020), which is based on the latest generation of Earth System Models which have participated in CMIP6 (Eyring et al., 2016). There is a potential inconsistency between our ECS, TCR and TCRE ranges, which arises because the ECS assessment comes from a study which uses multiple lines of evidence, the TCR assessment is based on a constrained set of CMIP6 models and the TCRE assessment is based on unconstrained CMIP6 Earth System Models. We discuss the importance of this inconsistency and its consequences in Section 4.

The other key metrics are related to effective radiative forcing (ERF, Forster et al., 2016). These values generally follow the AR5 assessment, except for aerosol, CO<sub>2</sub>, and methane ERF. For aerosol and CO<sub>2</sub> ERF, we use the more recent work of Smith, Kramer, Myhre, Alterskjær, et al. (2020). For methane ERF, we increase the AR5 assessment following Etminan et al. (2016) although we note that this increase may be offset by an updated understanding of the impact of rapid adjustments (Smith, Kramer, et al., 2018).

At this point, we stress that our proxy assessed ranges are only one of a range of possible choices. Assessing all the available literature is a demanding task that is well undertaken by the IPCC. We do not attempt to reproduce this task here. Instead, the key is that our proxy assessed ranges are (a) reasonable and (b) were available at the time of the study's inception.

Following this intercomparison consortium's choice of proxy assessed ranges, modeling groups then had the opportunity to develop parameter ensembles which best reflected these assessed ranges. As previously discussed, this allowed some modeling teams (although crucially not all) to use the same "constraining benchmarks" (with a number of different techniques being employed to consider the constraining benchmarks, see Table 1). We use these consistently constrained models to gain unique insights into the impact of differences in model structure and constraining techniques when RCMs are used as integrators of knowledge, free from a typical source of disagreement between the models, namely that they were constrained to reproduce different understandings of the climate. The inclusion of results from models which were not constrained using the same benchmarks allows us to quantify the importance of constraining when using reduced complexity climate models as integrators of knowledge.

The modeling groups submitted a range of concentration-driven, emission-driven and idealized scenarios for their chosen parameter subsets (see scenario specifics below). Subsequently, several metrics were calculated, such as TCR from the idealized CO<sub>2</sub>-only 1pctCO<sub>2</sub> experiment (in which atmospheric CO<sub>2</sub> concentrations rise at 1% per year from preindustrial levels). Calculating derived metrics on each individual ensemble member ensures that all metrics are calculated from internally self-consistent model runs, which is of particular importance when the metric is based on more than one output variable from the model (e.g., TCRE, which relies on both surface air temperature change and inverse emissions of CO<sub>2</sub>). If we instead calculated results based on percentiles of different variables, we would not be using an internally self-consistent set.

Where modeling groups felt it was more appropriate (e.g., OSCARv3.1), they performed their own weighting of ensemble members before submitting.

The one metric which is not easily calculated from model results is ECS because it is defined at equilibrium. Accordingly, modeling groups reported their own diagnosed ECS for each ensemble member, rather than performing experiments which would allow it to be calculated after submission had taken place.

When evaluating model performance, we are interested not only in how well a model can reproduce the best-estimate, but also the range, of a given quantity. A key part of any climate assessment is the uncertainty and it is critical that RCMs reflect the assessed likely and very likely ranges if they are to be used as integrators of knowledge. We assess the relative difference between the model and the assessed ranges at the very likely lower (fifth percentile, also referred to as “vll”), likely lower (17th percentile, “ll”), central (50th percentile, “c”), likely upper (83rd percentile, “lu”), and very likely upper (95th percentile, “vlu”). Assessing deviations using relative differences allows us to quickly evaluate how models perform over a range of metrics on the same scale.

The set of scenarios that each modeling group was asked to run follow the experimental protocols of CMIP6's ScenarioMIP (O'Neill et al., 2016). The SSPX-Y.Y experiments (e.g., SSP1-1.9, SSP2-4.5, SSP5-8.5) are defined in terms of concentrations of well-mixed greenhouse gases, i.e., CO<sub>2</sub>, CH<sub>4</sub>, N<sub>2</sub>O, hydrofluorocarbons (HFCs), perfluorocarbons (PFCs), and hydrochlorofluorocarbons (HCFCs), emissions of “aerosol precursor species emissions”, i.e., sulfur, nitrates, black carbon, organic carbon and ammonia and natural effective radiative forcing variations. As described in Z. R. J. Nicholls et al. (2020), where required, models may use prescribed effective radiative forcing if they do not include the required gas cycles or radiative forcing parameterizations.

The esm-SSPX-Y.Y experiments are identical to the SSPX-Y.Y experiments, except CO<sub>2</sub> emissions are prescribed instead of CO<sub>2</sub> concentrations, following the CMIP6 C4MIP protocol (Jones et al., 2016). Finally, we also perform esm-SSPX-Y.Y-allGHG experiments. These are identical to the esm-SSPX-Y.Y experiments, except they are defined in terms of emissions of all well-mixed greenhouse gases, not only CO<sub>2</sub>, rather than concentrations. There is no equivalent of these esm-SSPX-Y.Y-allGHG experiments in the CMIP6 protocol, however it is these experiments which are of most interest to WG3, given that WG3 focuses on scenarios defined in terms of emissions alone. We use the data sources described in Z. R. J. Nicholls et al. (2020) to specify the inputs for each of these scenarios. The input data set compilations, comprising emission, scenario and forcing data, as well as the protocols are archived with Zenodo (Z. Nicholls & Lewis, 2021)—and can contribute to scientific studies beyond this intercomparison as they largely reflect the CMIP6 experimental designs.

The protocol designed for this study requires that each RCM modeling group runs every probabilistic ensemble member once for each scenario and then submits their output for further analysis. With nine modeling groups participating, this intercomparison project compiled a database of results containing thousands of runs for each RCM, from which we can calculate different warming, effective radiative forcing or ocean heat uptake percentiles for a wide range of scenarios.

## 4. Results and Discussion

### 4.1. Fit to Assessed Ranges

The ability of RCMs to match the assessed ranges varies (Table 3, Figure 1, Table S1, and Figures S1–S9). In general, the RCMs capture the central assessed values better than the likely and very likely ranges. Historical warming, TCR and the TCRE are notable exceptions to this. For the TCR, the upper likely and very likely upper assessed values are captured by the RCMs about as well as the central value. For TCRE and historical warming, the very likely lower and likely lower assessed values are better captured by the RCMs than the central values.

Considering the variation between metrics, we see that the proxy assessment of the ocean heat content and effective radiative forcing metrics is better captured by the RCMs than the other metrics. For the ocean heat content and effective radiative forcing metrics, the median multimodel difference is less than or equal to 10% for the central proxy assessed range. However, there is less close agreement with the very likely and

**Table 3**  
*Comparison of Each Model's Probabilistic Distribution With the Proxy Assessment*

Climate model Assessed range	Multimodel median of magnitude of relative differences				
	vll	ll	c	lu	vlu
2000–2019 GMST rel. to 1961–1990	7%		11%		25%
Equilibrium climate sensitivity	<b>16%</b>	<b>15%</b>	<b>12%</b>	<b>14%</b>	<b>20%</b>
Transient climate response	38%	18%	7%	4%	7%
Transient climate response to emissions	<b>9%</b>	<b>11%</b>	<b>20%</b>	<b>19%</b>	<b>20%</b>
2014 CO <sub>2</sub> effective radiative forcing		<b>5%</b>	<b>5%</b>	<b>1%</b>	
2014 Aerosol effective radiative forcing		<b>14%</b>	<b>10%</b>	<b>19%</b>	
2018 Ocean heat content rel. to 1971		<b>1%</b>	<b>1%</b>	<b>16%</b>	
2011 CH <sub>4</sub> Effective radiative forcing		<b>4%</b>	<b>6%</b>	<b>18%</b>	
2011 N <sub>2</sub> O Effective radiative forcing		<b>11%</b>	<b>2%</b>	<b>10%</b>	
2011 F-Gases effective radiative forcing		<b>2%</b>	<b>3%</b>	<b>4%</b>	

Climate model Assessed range	CICERO-SCM					EMGC					FaIR1.6				
	vll	ll	c	lu	vlu	vll	ll	c	lu	vlu	vll	ll	c	lu	vlu
2000–2019 GMST rel. to 1961–1990	<b>6%</b>		<b>12%</b>		<b>17%</b>	–30%		–11%		35%	14%		23%		37%
Equilibrium climate sensitivity	9%	4%	–2%	–12%	–22%	–43%	–42%	–38%	–28%	–12%	–16%	–11%	–3%	11%	32%
Transient climate response	33%	14%	1%	–6%	–12%						43%	23%	10%	5%	8%
Transient climate response to emissions											<b>17%</b>	<b>–3%</b>	<b>–10%</b>	<b>–11%</b>	<b>–12%</b>
2014 CO <sub>2</sub> effective radiative forcing		<b>15%</b>	<b>8%</b>	<b>2%</b>			<b>12%</b>	<b>6%</b>	<b>–0%</b>			<b>–2%</b>	<b>5%</b>	<b>14%</b>	
2014 Aerosol effective radiative forcing		28%	37%	64%			<b>16%</b>	<b>16%</b>	<b>16%</b>			<b>2%</b>	<b>0%</b>	<b>–4%</b>	
2018 Ocean heat content rel. to 1971		<b>0%</b>	<b>–0%</b>	<b>–0%</b>			–9%	2%	26%			–0%	12%	24%	
2011 CH <sub>4</sub> effective radiative forcing		12%	–12%	–27%			23%	–3%	–20%			<b>3%</b>	<b>–8%</b>	<b>–15%</b>	
2011 N <sub>2</sub> O effective radiative forcing		<b>13%</b>	<b>–6%</b>	<b>–20%</b>			25%	4%	–12%			<b>8%</b>	<b>–2%</b>	<b>–9%</b>	
2011 F-Gases Effective Radiative Forcing												<b>–1%</b>	<b>–3%</b>	<b>–4%</b>	

Climate model Assessed range	FaIRv2.0.0-alpha					Hector					MAGICCv7.5.1				
	vll	ll	c	lu	vlu	vll	ll	c	lu	vlu	vll	ll	c	lu	vlu
2000–2019 GMST rel. to 1961–1990	11%		22%		38%	7%		16%		25%	<b>1%</b>		<b>1%</b>		<b>2%</b>
Equilibrium climate sensitivity	–22%	–15%	–2%	15%	31%	<b>–20%</b>	<b>–17%</b>	<b>–8%</b>	<b>0%</b>	<b>16%</b>	<b>–13%</b>	<b>–15%</b>	<b>–14%</b>	<b>–16%</b>	<b>–17%</b>
Transient climate response	26%	12%	4%	3%	4%	45%	25%	11%	3%	0%	32%	16%	7%	2%	0%
Transient climate response to emissions	<b>–1%</b>	<b>–16%</b>	<b>–20%</b>	<b>–19%</b>	<b>–20%</b>						<b>9%</b>	<b>–2%</b>	<b>–1%</b>	<b>1%</b>	<b>–2%</b>
2014 CO <sub>2</sub> effective radiative forcing		<b>3%</b>	<b>8%</b>	<b>14%</b>								<b>–2%</b>	<b>–1%</b>	<b>–1%</b>	
2014 Aerosol effective radiative forcing		–12%	–13%	–27%			51%	44%	29%			<b>–1%</b>	<b>–7%</b>	<b>–19%</b>	
2018 Ocean heat content rel. to 1971												<b>1%</b>	<b>1%</b>	<b>1%</b>	
2011 CH <sub>4</sub> effective radiative forcing		<b>5%</b>	<b>–1%</b>	<b>–5%</b>								<b>–1%</b>	<b>–3%</b>	<b>–3%</b>	
2011 N <sub>2</sub> O Effective radiative forcing		<b>4%</b>	<b>–2%</b>	<b>–7%</b>								<b>1%</b>	<b>2%</b>	<b>2%</b>	
2011 F-Gases effective radiative forcing		<b>3%</b>	<b>5%</b>	<b>7%</b>								<b>1%</b>	<b>1%</b>	<b>1%</b>	

**Table 3**  
Continued

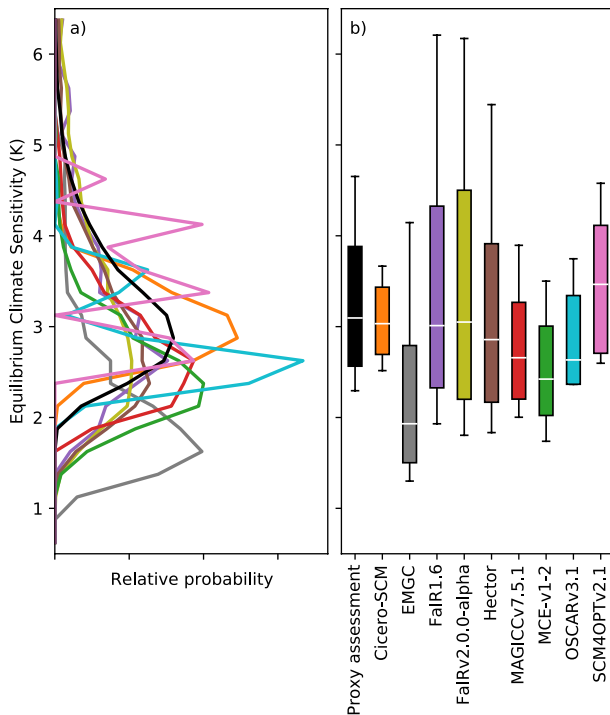
Climate model Assessed range	MCE-v1-2					OSCARv3.1					SCM4OPTv2.1				
	vll	ll	c	lu	vlu	vll	ll	c	lu	vlu	vll	ll	c	lu	vlu
2000–2019 GMST rel. to 1961–1990	<b>−1%</b>		<b>3%</b>		<b>3%</b>	<b>3%</b>		<b>2%</b>		<b>0%</b>	−8%		10%		28%
Equilibrium climate sensitivity	−24%	−22%	−22%	−23%	−26%	<b>3%</b>	<b>−9%</b>	<b>−15%</b>	<b>−14%</b>	<b>−20%</b>	<b>13%</b>	<b>4%</b>	<b>12%</b>	<b>6%</b>	<b>−3%</b>
Transient climate response	23%	4%	−6%	−12%	−15%	44%	21%	−1%	−10%	−15%	61%	35%	9%	0%	−6%
Transient climate response to emissions	0%	−18%	−23%	−23%	−27%	11%	−11%	−21%	−24%	−28%					
2014 CO <sub>2</sub> Effective radiative forcing		<b>−1%</b>	<b>−0%</b>	<b>1%</b>			<b>13%</b>	<b>6%</b>	<b>0%</b>			<b>6%</b>	<b>−0%</b>	<b>−6%</b>	
2014 Aerosol effective radiative forcing		<b>10%</b>	<b>6%</b>	<b>−9%</b>			−21%	−10%	−1%		−14%	−10%	−23%		
2018 Ocean heat content rel. to 1971		<b>−1%</b>	<b>−0%</b>	<b>1%</b>			−20%	5%	28%		−21%	0%	16%		
2011 CH <sub>4</sub> Effective radiative forcing		<b>−1%</b>	<b>−3%</b>	<b>−4%</b>			5%	−17%	−32%		2%	−19%	−34%		
2011 N <sub>2</sub> O Effective radiative forcing		<b>−1%</b>	<b>−0%</b>	<b>−1%</b>			26%	5%	−11%		21%	0%	−14%		
2011 F-Gases effective radiative forcing		<b>0%</b>	<b>−0%</b>	<b>0%</b>			<b>15%</b>	<b>4%</b>	<b>−6%</b>		27%	14%	4%		

*Note.* In each square, we show the relative difference between the model result and the proxy assessed value ( $\Delta_m$ , calculated as  $\Delta_m = \frac{m - a}{|a|}$  where  $m$  is the value from the model's probabilistic distribution and  $a$  is the proxy assessment value). Bold cells indicate that this model is within 20% of the proxy assessment at all likelihood levels for this metric. If a row is completely empty for a model, this indicates that the model did not submit results which allowed that metric to be calculated. Empty cells within a row which is otherwise not completely empty for a model indicates that no proxy assessment at this likelihood level was available (e.g., we have proxy assessments for likely lower 2014 CO<sub>2</sub> effective radiative forcing, but not for very likely lower 2014 CO<sub>2</sub> effective radiative forcing). Only the magnitude of  $\Delta_m$  from each model was used to calculate the multimodel median (to ensure that positive and negative values of  $\Delta_m$  from different models would not cancel out). The assessed ranges are labeled as “vll” (very likely lower, i.e., fifth percentile), “ll” (likely lower, 17th percentile), “c” (central, 50th Percentile), “lu” (likely upper, 83rd percentile) and “vlu” (very likely upper, 95th percentile).

likely proxy assessed ranges for the effective radiative forcing metrics, with median multimodel differences being up to 19% (aerosol effective radiative forcing).

For the other metrics (historical warming, ECS, TCR, and TCRE), the median multimodel difference is greater than 20% for at least one of the assessed ranges. However, there is significant variation across the likelihood levels. For example, the multimodel median matches the very likely lower historical warming (rows labeled “2000–2019 GMST rel. to 1961–1990” in Table 3) to within 7%. However, the multimodel median differs from the central and very likely upper historical warming by 11% and 25%, indicating that the models are having greater difficulty capturing the upper-end warming estimates.

There is also significant spread in performance across the models. MAGICCv7.5.1 performs better than the multimodel median across all metrics and assessed ranges (very likely lower, likely lower, central, likely upper, very likely upper) except for ECS while MCE-v1-2 performs better than the multimodel median across all metrics and assessed ranges except for three metrics (ECS, TCR and TCRE). However, all RCMs had at least one metric where they matched the proxy assessment at all likelihood levels to within 20% (bold cells in Table 3). For many applications, agreement to within 20% will be sufficient given the uncertainty associated with assessed ranges. However, for some applications, using an RCM's probabilistic distribution which has differences greater than 5–10% (for certain metrics) may be problematic as such differences could bias projections to an unacceptably large degree. For example, the WG3 classification of scenarios in terms of their peak warming levels should ideally be consistent with the range of evidence assessed in IPCC WG1. To have confidence that such an application is reflecting the WG1 assessment, the RCMs should be within 5–10% of the assessed results (particularly for any future warming assessment).



**Figure 1.** Distribution of equilibrium climate sensitivity (ECS) from each RCM (colored lines) and the proxy assessed range (solid black line). (a) Distribution of ECS; (b) very likely (whiskers), likely (box), and central (white solid line) from the proxy assessment and each RCM.

When interpreting these results it is vital to keep in mind that, for some models, the same benchmarks are used to both constrain and evaluate the models. The reason for this choice is that we are evaluating the ability of the models to act as integrators of knowledge, i.e., to simultaneously capture all the independent assessments (see also discussion in Section 1). We are not attempting to do a calibration followed by an out-of-sample evaluation.

As a result, it is not so surprising that the models which calibrated to the benchmarks, specifically MAGICC7 and MCE-v1-2, better reflect the benchmarks during the evaluation phase. However, the results presented here highlight just how important it is to calibrate if the model is to be used as an integrator of knowledge. If the goal is an integrator of knowledge which reflects key benchmarks, our results suggest that models which are calibrated will perform better.

## 4.2. Projection Results

For each probabilistic setup, the RCMs also submitted projections of global-mean surface temperature, effective radiative forcing (split into total, aerosols and CO<sub>2</sub>) and atmospheric CO<sub>2</sub> concentrations for the SSPX-Y.Y, ESM-SSPX-Y.Y, and ESM-SSPX-Y.Y-allGHG experiments.

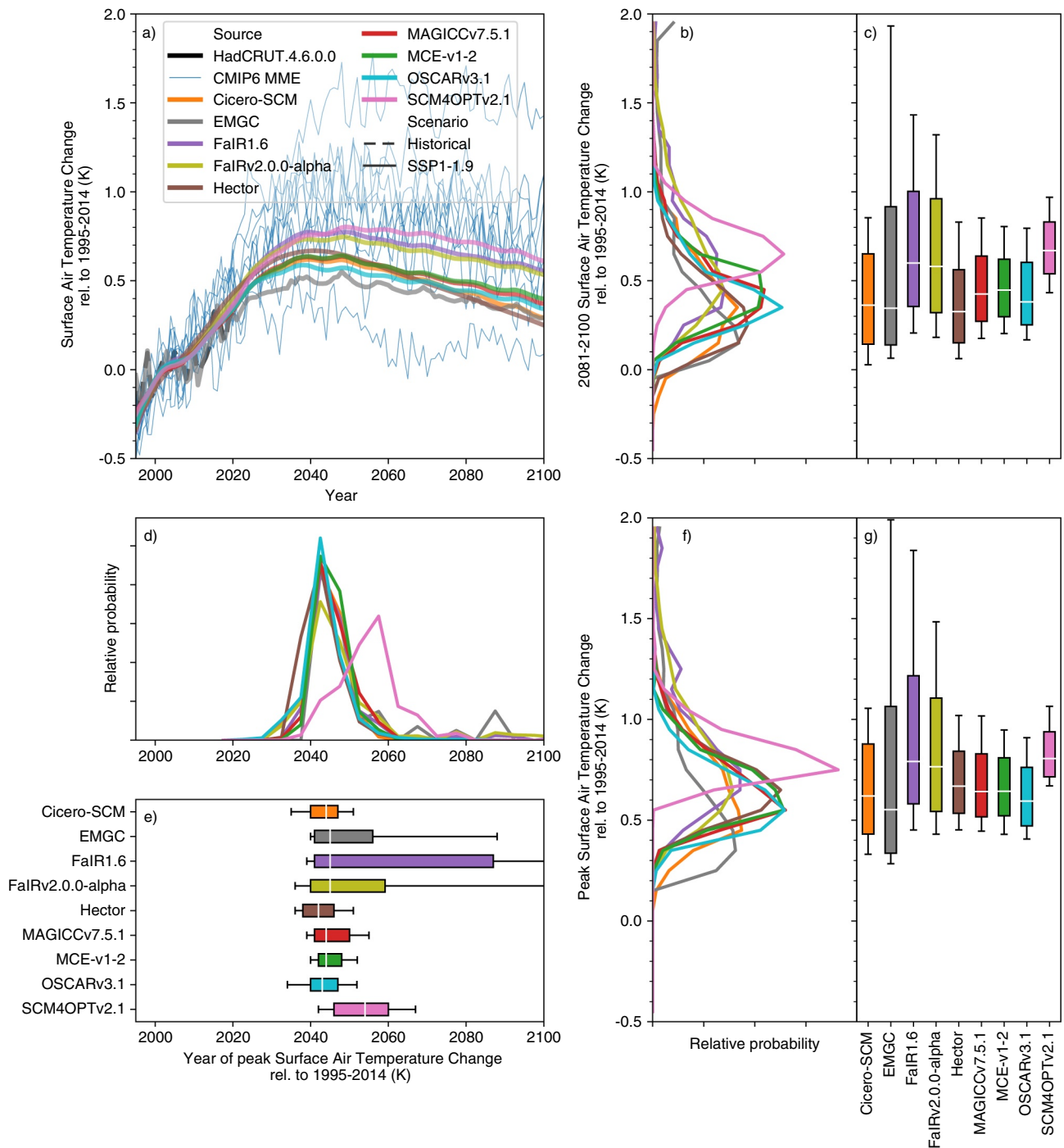
### 4.2.1. Global-Mean Surface Air Temperature

Under SSP1-1.9, median end of century (2081–2100) projections relative to 1995–2014 vary by 0.4°C across the models (from Hector with 0.3°C of warming to SCM4OPTv2.1 with 0.7°C, Figures 2a–2c). Variations in fifth percentile warming show a similar range, from 0.0 to 0.4°C. In contrast, upper-end, 95th percentile warming shows far greater variation, from 0.8 to 1.9°C. For the SSP1-1.9 scenario, the spread in RCMs' probabilistic projections is similar to the spread in the CMIP6 multimodel ensemble. Nonetheless, the most extreme CMIP6 model projections are outside the range of most RCMs' 5–95th percentiles. We discuss reasons for this difference in Section 4.3.

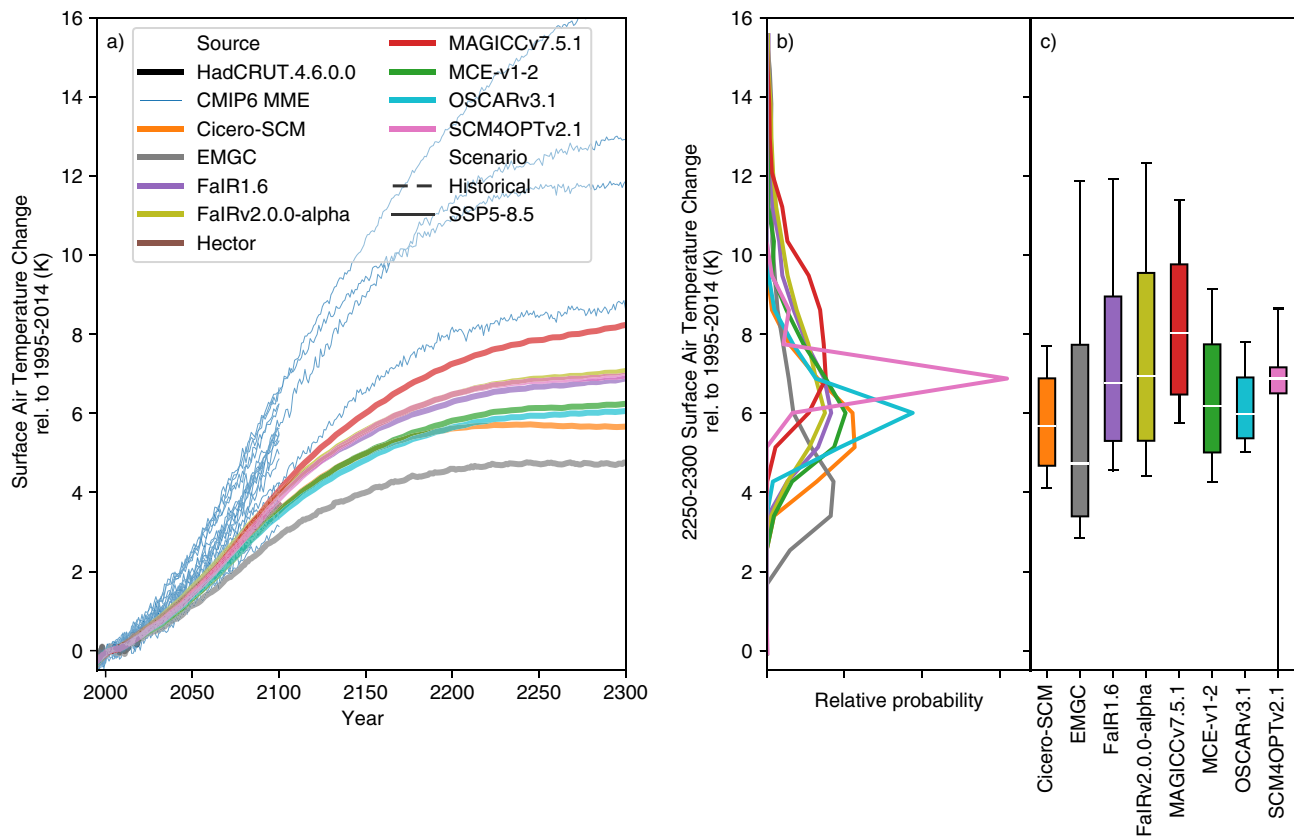
A slightly smaller spread is seen in peak temperature (Figures 2f and 2g). Across the RCM ensemble, SSP1-1.9 median peak warming ranges from 0.55 to 0.8°C while the 5th and 95th percentiles range from 0.3 to 0.7°C and 0.9 to 2.0°C, respectively. The year of peak warming shows much more variation, particularly at the upper end (Figures 2d and 2e). While the median peak year is fairly consistent across the RCMs' ensembles, around 2045 (although SCM4OPTv2.1's 2055 peak is a clear outlier), and the fifth percentile peak year varies from 2030 to 2040, the 95th percentile varies from 2050 to beyond the end of this century. In the EMGC, FaIR1.6 and FaIRv2.0.0-alpha probabilistic distributions, there is a significant area of parameter space which results in ongoing warming even after CO<sub>2</sub> emissions have reached net zero. These models also drive the spread in end of century temperature projections, particularly in the 95th percentile (Figures 2b and 2c).

In the SSP1-2.6 scenario (Supplementary Figure S10), median peak warming ranges from 0.65 to 1.1°C (0.1–0.3°C higher than in SSP1-1.9). Median end of century warming (relative to 1995–2014) ranges from 0.6 to 1.0°C. End of century fifth percentile warming ranges from 0.2 to 0.8°C and 95th percentile warming ranges from 1.2 to 2.0°C. As in SSP1-1.9, a number of CMIP6 model projections lie above the upper end of the constrained RCMs.

Under SSP1-2.6, the RCMs diverge more in their peak temperature projections, both compared to end of century warming and compared to SSP1-1.9. Once again, the fifth percentile and median are fairly consistent (ranging from 0.3 to 0.9°C and 0.65 to 1.1°C, respectively). However, 95th percentile projections vary from 1.2 to 2.8°C. The divergence in upper-end warming between SSP1-2.6 and SSP1-1.9 is driven by FaIR1.6, and appears to be the result of persistent warming after CO<sub>2</sub> emissions reach net zero given that



**Figure 2.** Surface air temperature (also referred to as global-mean surface air temperature, GSAT) change under the very-low-emissions SSP1-1.9 scenario. (a) GSAT projections from 1995 to 2100. We show the median RCM projections (colored lines), GMST observations from HadCRUT4.6.0.0 (Morice, Kennedy, Rayner, & Jones, 2012) up to 2019 (dashed black line) and CMIP6 model projections (thin blue lines, we show a single ensemble member for each CMIP6 model to preserve the CMIP6 models' natural variability signal); (b) distribution of 2081–2100 mean GSAT from each RCM; (c) very likely (whiskers), likely (box) and central (white line) 2081–2100 mean GSAT estimate from each RCM; (d) as in (b) except for the year in which GSAT peaks; (e) as in (c) except for the year in which GSAT peaks; (f) as in (b) except for the peak GSAT; (g) as in (c) except for the peak GSAT. All results are shown relative to the 1995–2014 reference period.



**Figure 3.** Long-term surface air temperature (also referred to as global-mean surface air temperature, GSAT) change under the high-emissions SSP5-8.5 scenario. (a) GSAT projections from 1995 to 2300. We show the median RCM projections (colored lines), GMST observations from (Morice, Kennedy, Rayner, & Jones, 2012) up to 2019 (dashed black line) and available CMIP6 model projections (thin blue lines, we show a single ensemble member for each CMIP6 model to preserve the CMIP6 models' natural variability signal); (b) distribution of 2250–2300 mean GSAT from each RCM; (c) very likely (whiskers), likely (box), and central (white line) 2250–2300 mean GSAT estimate from each RCM. All results are shown relative to the 1995–2014 reference period.

its 83rd percentile peak warming year is after 2100. Across the models, peak warming year shows a similar range to SSP1-1.9, albeit occurring 25–30 years later in the median (ranging from 2065 to 2075). Once again, the fifth percentile (ranging from 2050 to 2060) shows a much smaller spread across the models than the 95th percentile (ranging from 2075 to beyond the end of the 21st Century).

The warmest RCMs in mitigation scenarios are also the warmest under the high-emissions, SSP5-8.5, scenario (Supplementary Figure S11). The exceptions to this are MAGICC7, which is one of the warmest models in SSP5-8.5 even though it was around the median in mitigation scenarios, and SCM4OPTv2.1, which was the warmest model in mitigation scenarios but is slightly cooler than the warmest models in SSP5-8.5. Under SSP5-8.5, median end of century warming ranges from 2.5 to 3.6°C across the RCMs. Unlike the mitigation scenarios, there is a similar level of disagreement in 5th and 95th percentile warming, with the fifth percentile ranging from 1.7 to 3.1°C and the 95th percentile ranging from 3.8 to 5.4°C. The RCMs all make future warming projections in the lower-half of the CMIP6 multimodel ensemble. Such a difference is largely explained by the constraints applied to the RCMs (see discussion in Section 4.3).

If we consider long-term (2250–2300) warming under the SSP5-8.5 scenario (Figure 3, see Supplementary Figures S12 and S13 for long-term warming under SSP1-1.9 and SSP1-2.6, respectively), the difference between RCMs and CMIP6 is even clearer (although the few CMIP6 models which have run the SSP5-8.5 extension are all at or above the median of the CMIP6 multimodel ensemble in 2100). On these timescales, MAGICC7 is clearly the warmest model, despite having slightly lower long-term effective radiative forcing than FaIR1.6, FaIR-v2.0.0-alpha, and MCE-v1-2 (Supplementary Figure S14). There is a significant spread in long-term projections across the RCMs, with the median ranging from 4.5 to 8.0°C, fifth percentile from 3°C (ignoring SCM4OPTv2.1 as an outlier) to 5.8°C and 95th from 7.8 to 12.3°C. Even these upper end

projections are well below the highest CMIP6 projections, which reach over 16°C of global-mean warming (again, likely due to constraining, see discussion in Section 4.3). Across all the RCMs, only CICERO-SCM shows any sign of temperatures peaking by 2300 under such a high-emissions scenario.

#### 4.2.2. Effective Radiative Forcing

Compared to temperatures, there is less variance in end of century total effective radiative forcing projections (Figure 4, Supplementary Figure S15 and S16). This finding reinforces the understanding that the parameterization of the climate response to effective radiative forcing is a key driver of climate projection uncertainty.

In SSP1-1.9, 2081–2100 mean total effective radiative forcing varies from 2.2 to 2.6 W/m<sup>2</sup>. The fifth percentile ranges from 1.8 to 2.1 W/m<sup>2</sup> across the models (excluding CICERO-SCM which has an extremely narrow range). The spread is larger for the 95th percentile, which ranges from 2.4 to 3.2 W/m<sup>2</sup>. This pattern, of uncertainty being higher for upper percentiles than lower percentiles, is seen across other key scenarios and highlights that the high-end effective radiative forcing projections are much more uncertain than the median and low-end effective radiative forcing projections.

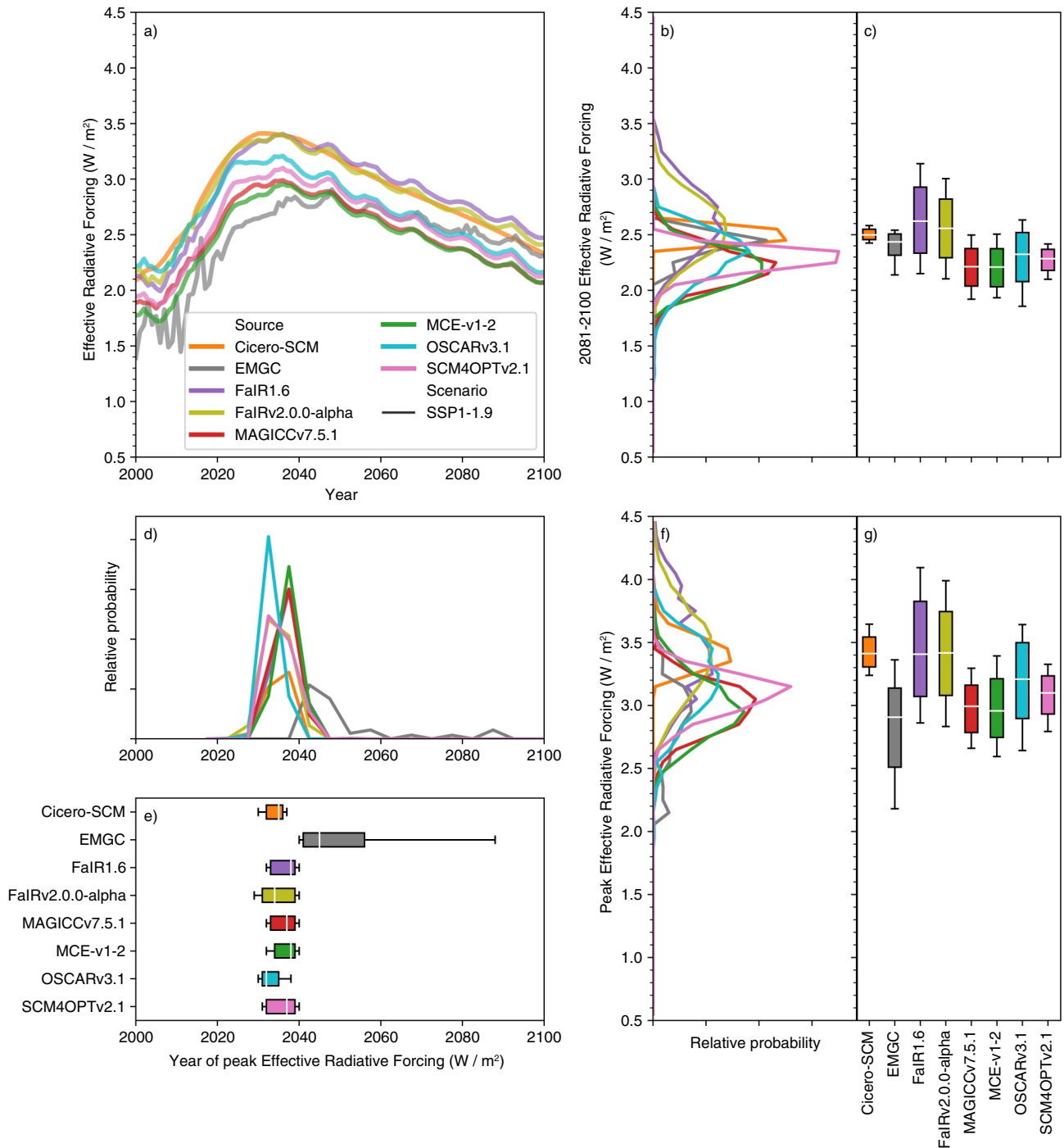
In SSP1-2.6 (Supplementary Figure S15, once again excluding CICERO-SCM because of its narrow range) median 2081–2100 total effective radiative forcing ranges from 2.9 to 3.4 W/m<sup>2</sup> while the fifth percentile only ranges from 2.4 to 2.7 W/m<sup>2</sup> and the 95th percentile has a much wider range of 3.1 to 4.1 W/m<sup>2</sup>. Under SSP5-8.5 (Supplementary Figure S16, excluding EMGC and CICERO-SCM as outliers), median 2081–2100 total effective radiative forcing ranges from 8.0 to 9.3 W/m<sup>2</sup> while the fifth percentile only ranges from 7.4 to 7.8 W/m<sup>2</sup> and the 95th percentile has a much wider range of 8.4 to 11.0 W/m<sup>2</sup>.

The approximate agreement in total effective radiative forcing is reflected in the agreement of each of the key contributors to this total, namely CO<sub>2</sub> and aerosol effective radiative forcing (Figure 5 and Supplementary Figures S17–S29, which also show ERF output up to the year 2300). The key exceptions to this are SCM4OPTv2.1 and OSCARv3.1's aerosol effective radiative forcing. This negative aerosol forcing is driven by SCM4OPTv2.1 and OSCARv3.1's inclusion of a climate feedback on aerosol effective radiative forcing. The climate feedback makes their median end of century aerosol effective radiative forcing 0.3–0.8 W/m<sup>2</sup> more negative than other RCMs across the scenarios, although the effect is stronger in OSCARv3.1 than in SCM4OPTv2.1. The strong aerosol forcing is somewhat compensated by other forcing agents although both these models have long-term ERF which is at the low end of the RCM ensemble under SSP5-8.5 (Supplementary Figure S14). The different aerosol ERF parameterizations warrant further attention, particularly because models without this aerosol ERF—climate feedback may be underestimating the spread in future temperature projections.

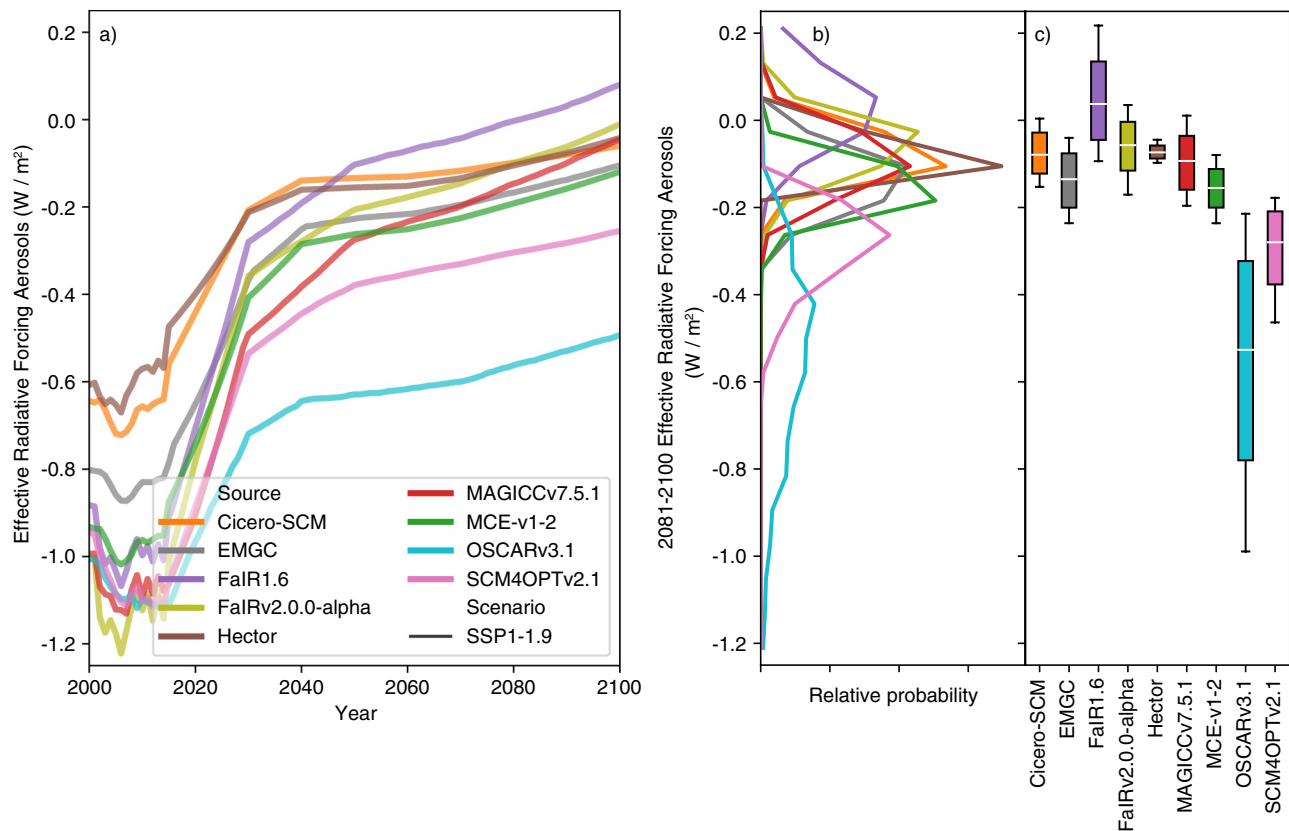
#### 4.2.3. Carbon Cycle

Moving beyond effective radiative forcing and its temperature response, we consider the behavior of the carbon cycle in the different RCMs. Clearly, the analysis presented here covers only a limited subset of the full range of carbon cycle behavior and metrics. The analysis is intended to highlight variance in carbon cycle behavior across the RCMs, providing the motivation for a more detailed future analysis. We use the emissions-driven ESM-SSPX-Y.Y set of scenarios, in which emissions of CO<sub>2</sub> are prescribed and atmospheric CO<sub>2</sub> concentrations are allowed to freely evolve (in contrast to the SSP experiments in which CO<sub>2</sub> concentrations are prescribed).

There are considerable variations between the RCMs which submitted relevant results (Supplementary Figures S30 and S31, Figure 6). In esm-SSP1-1.9 (Supplementary Figure S30, excluding CICERO-SCM because of its narrow range), the spread in median peak atmospheric CO<sub>2</sub> concentrations (430–450 ppm) is similar to the spread in 2081–2100 median concentrations (385–410 ppm). Similarly, in esm-SSP1-2.6 (Figure S31, again excluding CICERO-SCM), the spread in median peak atmospheric CO<sub>2</sub> concentrations (450–480 ppm) shows a spread similar to the spread in 2081–2100 median concentrations (430–460 ppm). Under both scenarios, there are wide variances in percentile ranges across the models, with MAGICC7 showing the largest uncertainty in 2081–2100 atmospheric CO<sub>2</sub> concentrations and SCM4OPTv2.1 showing the least (arguably, this model's range is overly confident). The considerable spread in projections from the models highlights the importance of carbon cycle uncertainty for emissions-driven projections. The spread



**Figure 4.** Effective radiative forcing under the very-low-emissions SSP1-1.9 scenario. (a) Median effective radiative forcing projections from 1995 to 2100 for each RCM; (b) distribution of 2081–2100 mean effective radiative forcing from each RCM; (c) very likely (whiskers), likely (box), and central (white line) 2081–2100 mean effective radiative forcing estimate from each RCM; (d) as in (b) except for the year in which effective radiative forcing peaks; (e) as in (c) except for the year in which effective radiative forcing peaks; (f) as in (b) except for the peak effective radiative forcing; (g) as in (c) except for the peak effective radiative forcing.



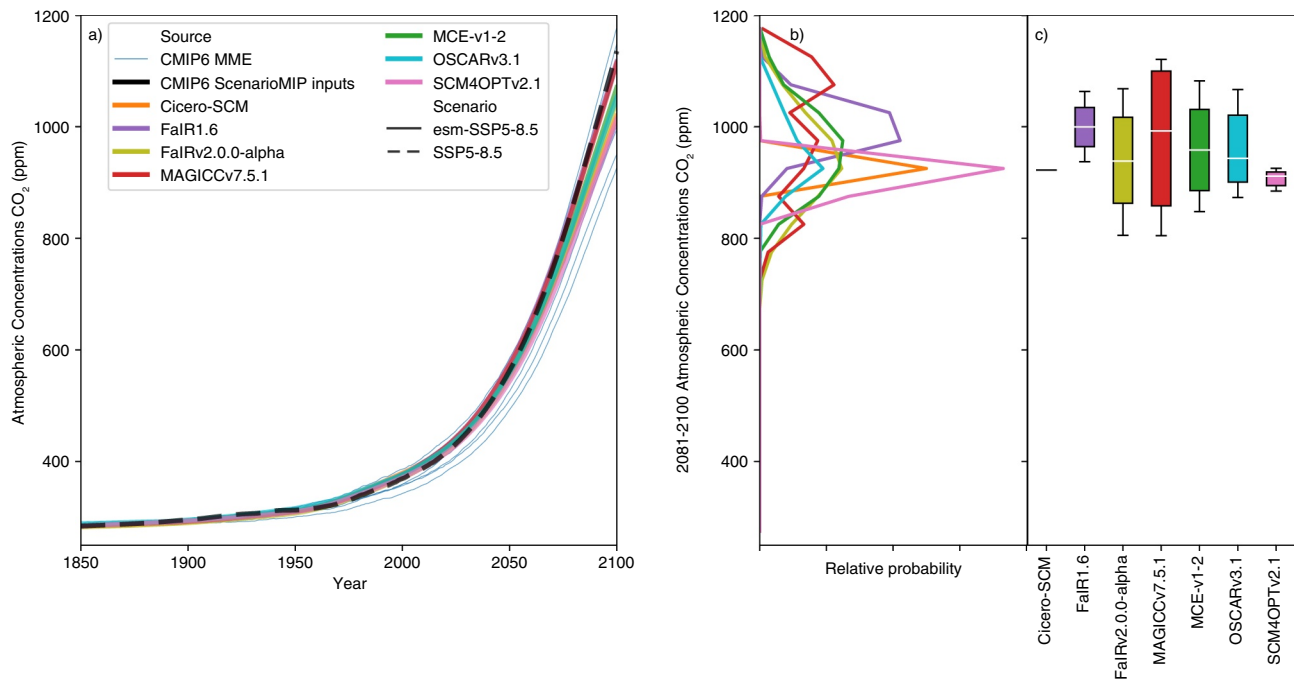
**Figure 5.** As in panels (a–c) of Figure 4, except for effective radiative forcing due to aerosols.

reinforces the need for a detailed study into available techniques for evaluating and potentially constraining carbon cycle behavior. Such a study would provide information about whether any of these projections can be ruled out based on other lines of evidence.

Next, we consider esm-SSP5-8.5, the only scenario with available CMIP6 Earth System Model results (Figure 6). Median 2081–2100 atmospheric CO<sub>2</sub> concentrations range from 920 ppm to 1,000 ppm while fifth percentile and 95th percentile concentrations range from 800 to 930 ppm and 910 ppm to 1,130 ppm respectively. MAGICC7 once again shows the largest uncertainties, but is more similar to the other RCMs than in the other scenarios. These comparisons highlight differences in the dynamics of the carbon cycle (and its feedbacks) in the various RCMs: uncertainties widen to a greater extent in higher-warming scenarios in FaIR1.6, FaIRv2.0.0-alpha, MCE-v1-2, OSCARv3.1, and SCM4OPTv2.1 compared to MAGICC7.

Median atmospheric CO<sub>2</sub> projections from all of the RCMs lie within the plume of available CMIP6 results (Figure 6). FaIR1.6 lies at the top end of the CMIP6 plume, and its 5–95th range does not include low-end CMIP6 results. In contrast, SCM4OPTv2.1 lies at the bottom end of the CMIP6 plume. FaIR-v2.0.0-alpha, MAGICC7, MCE-v1-2 and OSCARv3.1 approximately span the CMIP6 range, with FaIR-v2.0.0-alpha's and MCE-v1-2's ranges being almost exactly in line with the CMIP6 range while MAGICC7's projections are slightly wider than the CMIP6 range and OSCARv3.1's projections are slightly narrower than the CMIP6 range. CICERO-SCM does not include uncertainty in the carbon cycle, nor temperature feedbacks on the carbon cycle, hence produces only a single best-estimate projection.

Despite the limits of our carbon cycle evaluation, it is notable that the CMIP6 ScenarioMIP input concentrations are generally higher than the RCMs' medians in emissions-driven runs across all considered scenarios. Emissions-driven scenario data from CMIP6 ESMs is almost exclusively related to the esm-SSP5-8.5 experiment. Hence, while the pattern appears to be that the prescribed SSP5-8.5 CMIP6 concentrations are at the high-end of the range compared to the esm-SSP5-8.5 CMIP6 ESM results, there is little data with which to determine whether the prescribed CO<sub>2</sub> concentrations in the low-emissions scenarios would be within the

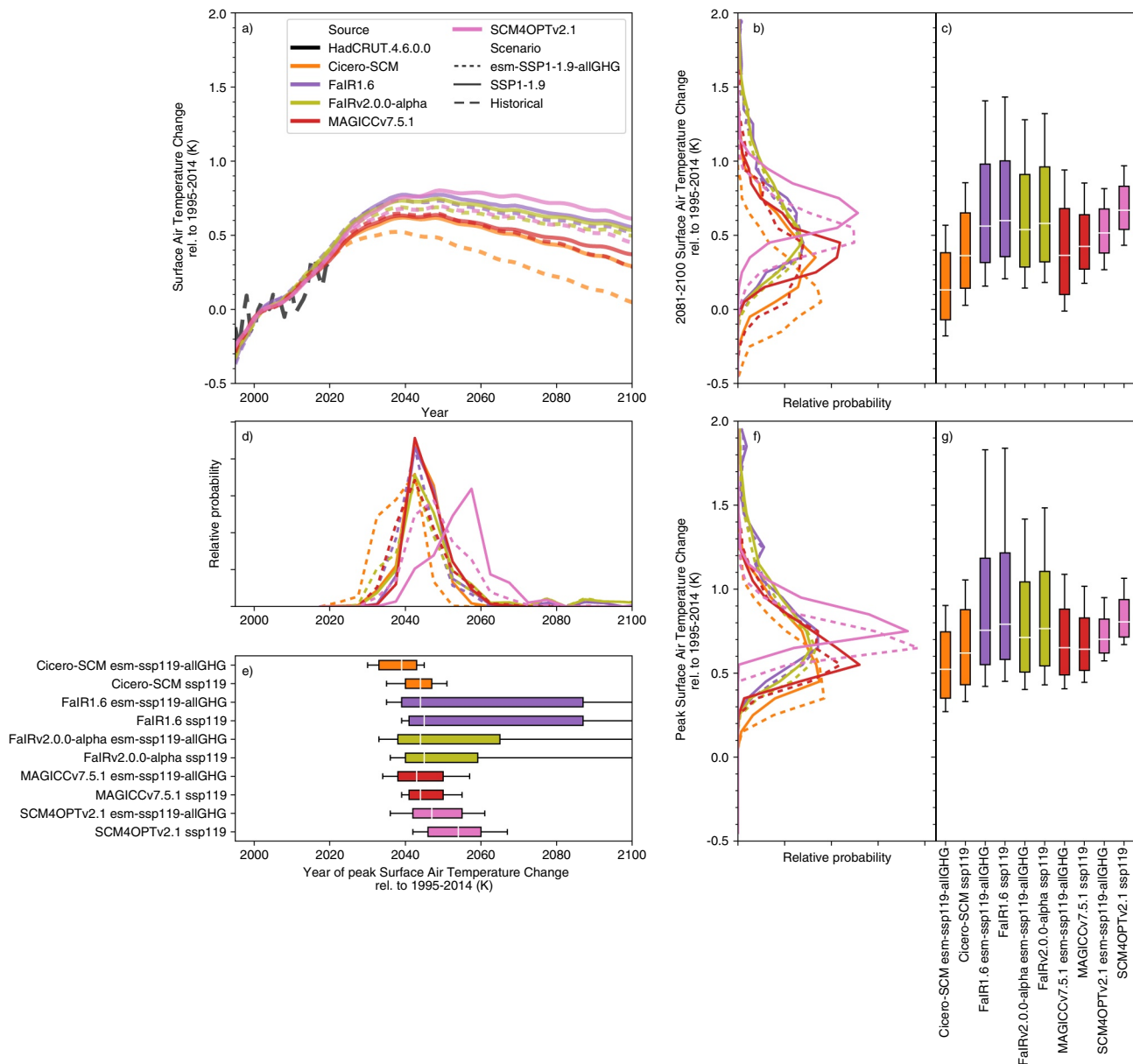


**Figure 6.** Atmospheric CO<sub>2</sub> concentration projections in the esm-SSP5-8.5 experiment. (a) Atmospheric CO<sub>2</sub> concentration projections from 1995 to 2100. We show the median RCM projections (colored lines), prescribed CMIP6 ScenarioMIP input concentrations from the SSP5-8.5 concentration-driven experiment (dashed black line) and available CMIP6 model projections (thin blue lines, we show a single ensemble member for each CMIP6 model to preserve the CMIP6 models' natural variability signal); (b) distribution of 2081–2100 mean atmospheric CO<sub>2</sub> concentration projections from each RCM; (c) very likely (whiskers), likely (box), and central (white line) 2081–2100 mean atmospheric CO<sub>2</sub> concentration projections estimate from each RCM. *Note.* that FaIR1.6 data is taken from the esm-SSP5-8.5-allGHG simulations because esm-SSP5-8.5 simulations are not available.

projected concentration change by emission-driven ESM models. In hindsight, the input atmospheric CO<sub>2</sub> concentrations used in the concentration-driven runs may turn out to be at the high-end of CMIP6 ESM results across a range of scenarios. Given that only one set of input concentrations can be used in CMIP6, it is not surprising that the CO<sub>2</sub> concentrations prescribed for CMIP6 experiments do not sit exactly in the middle of later emissions-driven runs. The opposite was observed in CMIP5: the input CO<sub>2</sub> concentrations (derived with MAGICC6) were found to be in the lower-half of the results from the CMIP5 emissions-driven runs (Friedlingstein et al., 2014). The CMIP6 concentrations were derived using an alpha version of MAGICC7, calibrated to approximately the median of the CMIP5 ESM carbon cycle responses with the inclusion of permafrost CO<sub>2</sub> and methane feedbacks (Meinshausen, Nicholls, et al., 2020). Choosing a carbon cycle parameterization more in line with the median of CMIP5 models appears to have led to CO<sub>2</sub> concentrations which are now in the upper-half of CMIP6 ESM projections (Figure 6). Whenever a single estimate of the relationship between CO<sub>2</sub> emissions and concentrations is used, there is always the risk that it will not be the central estimate of the next generation of ESMs as our understanding of the carbon cycle improves and the ensembles of participating ESMs changes in each intercomparison phase. While this does not invalidate the design of concentration-driven experiments which are developed in this way, it must be kept in mind when relating emissions scenarios and the output of concentration-driven CMIP experiments.

#### 4.2.4. All Greenhouse Gas Emissions-Driven Runs

The final set of experiments we present are the experiments which are most relevant to WG3: all greenhouse gas emissions-driven runs. As discussed in Section 1, WG3 describes scenarios in terms of their emissions hence needs models which can run in a fully emissions-driven setup. The cost of running ESMs for a large number of scenarios and parameter configurations in such a setup is computationally prohibitive (and few ESMs include key feedbacks such as methane permafrost and wetland emissions), hence there is a paucity of data against which to evaluate the projections of RCMs in such experiments. Nonetheless, here we pres-



**Figure 7.** Surface air temperature (also referred to as global-mean surface air temperature, GSAT) change in the concentration-driven SSP1-1.9 experiment and the all greenhouse gas emissions-driven esm-SSP1-1.9-allGHG experiment. (a) GSAT projections from 1995 to 2100. We show the median RCM projections (colored lines) for the concentration-driven experiment (solid) and all greenhouse gas emissions-driven experiment (dashed) as well as observations up to 2019 (dashed black line); (b) distribution of 2081–2100 mean GSAT for each scenario from each RCM; (c) very likely (whiskers), likely (box), and central (white line) 2081–2100 mean GSAT estimate for each scenario from each RCM; (d) as in (b) except for the year in which GSAT peaks; (e) as in (c) except for the year in which GSAT peaks; (f) as in (b) except for the peak GSAT; (g) as in (c) except for the peak GSAT. All results are shown relative to the 1995–2014 reference period.

ent the results of such experiments in the hope that they will inspire further efforts into how to validate RCMs in this fully coupled, all greenhouse gas emissions-driven setup.

Five models (CICERO-SCM, FaIR1.6, FaIRv2.0.0-alpha, MAGICC7, and SCM4OPTv2.1) have submitted results for the all greenhouse gas emissions-driven experiments. The results suggest that the all greenhouse gas emissions-driven runs are cooler and peak earlier than the concentration-driven runs (Figure 7, Supplementary Figure S32 and S33). However, the magnitude of the difference varies across the models. For median projections, MAGICC7 suggests the smallest difference between concentration-driven and all green-

house gas emissions-driven runs while CICERO-SCM and SCM4OPTv2.1 imply differences of up to 0.3°C for peak and 2081–2100 warming and a peak in warming up to 10 years earlier. The range of projections in the all greenhouse gas emissions-driven runs are generally about the same or slightly wider than in the concentration-driven runs, with MAGICC7 showing the largest increase in projection ranges.

The lower-warming and wider projection ranges seen in all greenhouse gas emissions-driven runs are consistent with two other pieces of knowledge. The first is that median CO<sub>2</sub> concentrations are lower in all greenhouse gas emissions-driven runs than in concentration-driven runs (Section 4.2.3). The second is that carbon cycle and other greenhouse gas cycle uncertainties are included in temperature projections in all greenhouse gas emissions-driven runs, while these uncertainties are missing in concentration-driven runs. The difference between the all greenhouse gas emissions-driven runs and concentration-driven runs reinforces the need for further consideration of RCM behavior beyond the climate response to ERF.

### 4.3. Further Discussion

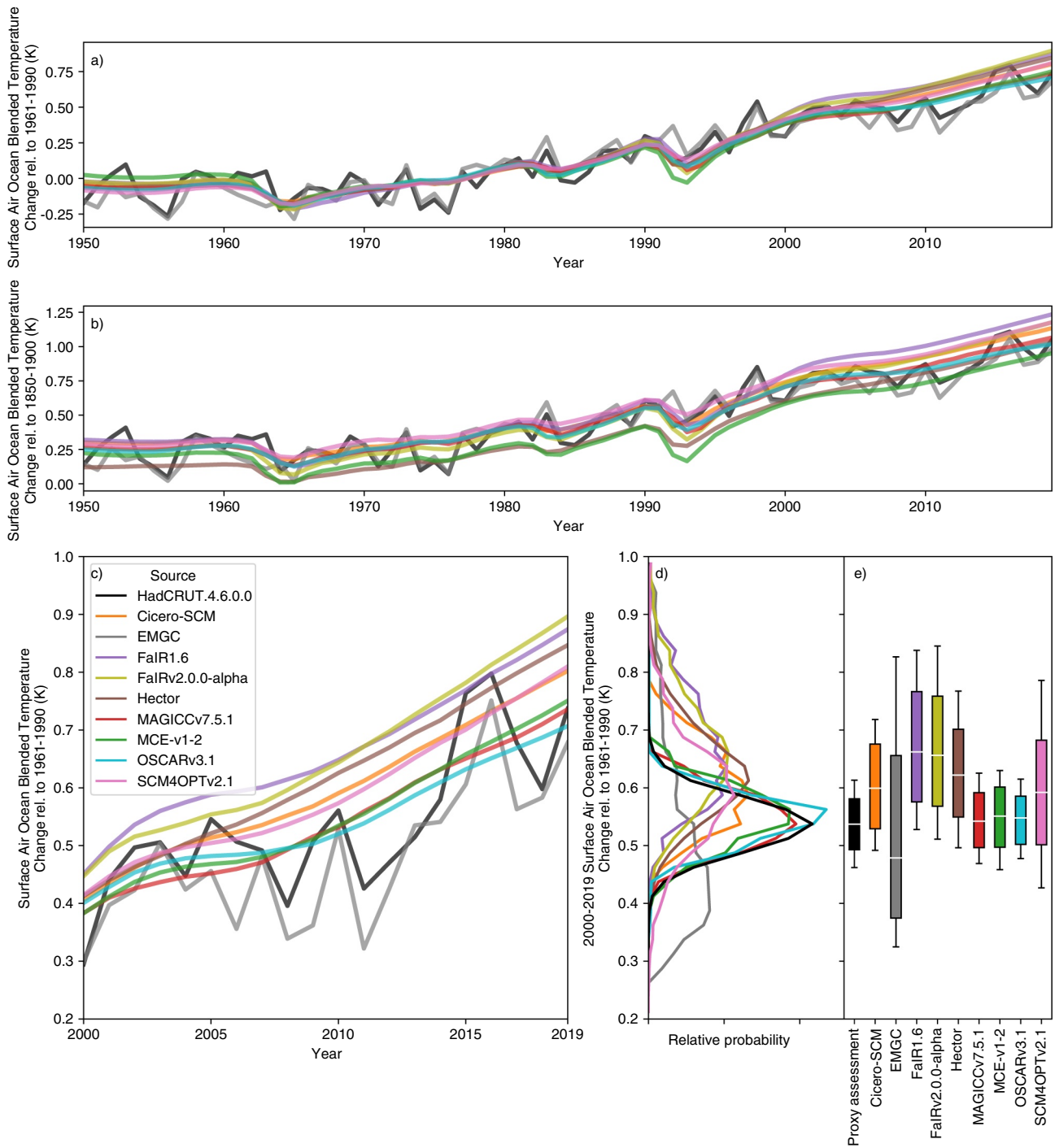
Our results prompt consideration of a number of further points. First, the assessment performed here provides a way to easily identify differences between an RCM's behavior and the assessed range of a particular metric. Such differences are important to quantify, as they can reveal biases in a probabilistic distribution. The quantification makes it possible for the users of these distributions to identify where the biases might impact their own conclusions.

There are, however, cases where the issue lies in the combination of the proxy assessed ranges taken together, rather than in the probabilistic distributions. In this study, we used a combination of ECS from the literature (based on multiple lines of evidence), TCR from constrained CMIP6 models and TCRE from unconstrained CMIP6 Earth System Models. This combination is likely to be slightly inconsistent. Unfortunately, inconsistency between metric values is an inevitable risk of using independent lines of evidence. The potential inconsistency could in part explain our finding that the RCMs' TCR ranges are generally too high, while their ECS and TCRE ranges are generally too low. To explain the inconsistency in more detail, first consider the ratio between TCR and ECS, i.e., the realized warming fraction. The realized warming fraction implied by our TCR and ECS distributions is around 0.5. This is at the low end of the assessment by Millar, Otto, et al. (2015). Hence, it can be argued that greater consistency within the proxy assessment would be achieved if either our proxy assessed TCR values were larger, or our proxy assessed ECS values were smaller. Similarly, the airborne fraction implied by our TCR and TCRE assessment is around 0.65. This is at the high-end of the CMIP5 and CMIP6 range quantified by Arora, Katavouta, et al. (2020). Once again, it can be argued that greater consistency within the proxy assessment would be achieved if either our proxy assessed TCR values were larger, or our proxy assessed TCRE values were smaller. Identifying such inconsistencies is a useful secondary benefit of exercises such as the one performed here.

Next, while they are a useful way of quickly visualizing a model's agreement with the (here proxy) assessed ranges, summary tables of the form of Table 3 hide the full story. Specifically, for timeseries based variables, assessed ranges can only consider the trend or change between specific reference periods and don't consider the entire timeseries as a whole.

Not considering the entire timeseries can lead to problematic interpretations of the agreement between a model and the assessment. A clear example here is historical surface air ocean blended temperature change. In our proxy assessment, we focused on 2000–2019 warming relative to the 1961–1990 reference period. On this measure, many of the RCMs were too warm compared to observations. However, the level of agreement is clearly reference period dependent (Figures 8a and 8b). In Figure 8a, which uses a 1961–1990 reference period, MAGICC7, MCE-v1-2, and OSCARv3.1 show the best agreement with observations (as also seen in Table 3). However, if we use a different reference period, e.g., 1850–1900 (Figure 8b), that impression changes with Hector, MAGICC7, and OSCARv3.1 being the closest to observations in the recent period.

Considering the entire timeseries provides a more robust check on model behavior. Fitting only to one evaluation and reference period can be achieved by slightly adjusting different model behavior, e.g., aerosol effective radiative forcing. However, if the entire timeseries are considered with multiple reference periods,



**Figure 8.** Historical surface air ocean blended temperature change (also referred to as global-mean surface temperature, GMST) from each RCM. We compare observations from HadCRUT4.6.0.0 (Morice, Kennedy, Rayner, & Jones, 2012) (solid black line) to the distribution from each RCM (colored lines). All panels use 1961–1990 as the reference period, the same reference period as is used in our proxy assessed ranges, except (b) which uses 1850–1900. (a, b) median GMST from 1950 to 2019; (c) median GMST from 2000 to 2019 (the proxy assessment period); (d) distribution of 2000–2019 mean GMST from each RCM and the proxy assessed range; (e) Very likely (whiskers), likely (box), and central (white line) estimate of 2000–2019 mean GMST from each RCM and the proxy assessed range. The historical simulation has been extended with SSP2-4.5 for the period 2015–2019.

such tuning quickly becomes impossible and the check provides detail into how well a model's dynamics are consistent with observations.

Moving away from evaluating the models, we find that higher historical warming, ECS and TCR values generally lead to higher-warming projections (an intuitive result). Hector provides an exception to this pattern, with relatively low temperature projections, especially in SSP1-1.9, despite its relatively high historical warming and TCR.

In the strong mitigation scenarios (SSP1-1.9 and SSP1-2.6), there is agreement to within  $\sim 0.1^{\circ}\text{C}$  in future projections (both best-estimate and range) between the models which best reflect historical warming (MAGICC7, MCE-v1-2, and OSCARv3.1). This agreement suggests that constraining greatly increases confidence in future projections. However, a limited set of models also provided probabilistic distributions that are constrained to match HadCRUT.5.0.1.0 (Morice, Kennedy, Rayner, Winn, et al., 2021), which is significantly warmer than the HadCRUT.4.6.0.0 based constraint used in the rest of the study. The future projections from these HadCRUT.5.0.1.0-constrained distributions are noticeably warmer (Supplementary Figures S34–S36) than projections from HadCRUT.4.6.0.0-constrained distributions, which demonstrates that projections are sensitive to the choice of constraint.

Given the sensitivity of conclusions to the constraint, the use of constraints must be carefully considered as it could lead to overconfidence (Sanderson et al., 2017). Even though considerable care is taken both here and elsewhere to identify and use relevant, physically justifiable, constraints, it is still possible that future research may show that the constraints are leading to overconfident future projections. Having said this, Herger et al. (2019) suggest that using multiple constraints, as is done by many RCMs here, reduces the likelihood of overconfidence.

Studies which constrain the raw CMIP6 model ensemble help explain the difference between the RCM-based results presented here and the raw CMIP6 model ensemble. Brunner et al. (2020), Liang et al. (2020), and Tokarska et al. (2020) all find significant reductions in both the best-estimate and 5–95% range GSAT projections after applying observed-warming constraints to the CMIP6 model ensemble. For the SSP1-2.6 and SS5-8.5 scenarios respectively, these studies find 5–95% GSAT (relative to 1995–2014) ranges of: Tokarska et al. (2020): 0.41–1.46 and 2.26–4.60 $^{\circ}\text{C}$ ; Liang et al. (2020) 0.52–1.66 and 2.72–4.77 $^{\circ}\text{C}$  and Brunner et al. (2020) 0.61–1.85 and 2.72–4.86 $^{\circ}\text{C}$ . These estimates, particularly for the SSP1-2.6 scenario, are slightly wider than our results based on RCMs. However, the constrained CMIP6 estimates are much closer to our RCM-based estimates than the raw CMIP6 model ensemble, in particular for the 95th percentile. This suggests that the majority of the difference between our RCM-based results and the raw CMIP6 model ensemble is explained by the constraining applied to the RCMs, rather than structural differences between RCMs and CMIP6 models (although structural differences may explain the disagreement between constrained CMIP6 output and our results). Further studies are needed to explore the validity of the constraining approaches for both ESMs and RCMs—as investigated here—but this study lays the foundation for systematically investigating probabilistic RCM ensembles in more detail.

Given the proxy assessment and results, we make one final observation: to extrapolate assessed warming ranges from one set of scenarios (e.g., the RCP or SSP-based scenarios) to a wider set of scenarios, it may be beneficial to include a benchmark of assessed future warming under a selection of scenarios. This benchmark could be taken from other studies, e.g., those that constrain CMIP projections (for the limited number of scenarios run by CMIP) based on historical observations (e.g., Brunner et al., 2020; Liang et al., 2020; Tokarska et al., 2020). Adding such a benchmark to the historical observations, present-day assessments and idealized metrics used in this study would highlight where future warming significantly diverges from other lines of evidence. Including scenarios with similar end of century total ERF but different transient evolutions (like the SSP4-3.4 and SSP5-3.4-overshoot scenario pair) would provide an even stronger check of the models' transient response. Such quantifications could be key when assessing future projections under large sets of scenarios, like the WG3 scenario database climate assessment. Of course, the risk of adding such benchmarks is an artificial narrowing of uncertainty in projected warming. Hence, future projections should only be included where there is a clear need and justification for consistency between the RCMs' projections and the projections from other lines of evidence.

## 5. Future Work

This exercise is a first step toward more comprehensive, routine evaluation of RCMs' probabilistic parameter ensembles and their corresponding projections. However, there is still much room for future work to improve on this study and the first phase of RCMIP. As a first suggestion, repeating this exercise with the assessed ranges from Working Group 1 of the Intergovernmental Panel on Climate Change's Sixth Assessment Report (due in mid-2021) would provide an evaluation of the extent to which RCMs can capture the latest international assessment of the scientific literature.

This future work could go beyond evaluation and also diagnose the root causes of differences between the models. One obvious area for examination would be the aerosol ERF, particularly the inclusion of a climate feedback in aerosol ERF parameterizations. Such an exercise could also provide greater insights into differences between the constrained RCMs' probabilistic distributions, the raw CMIP6 multimodel ensemble and constrained CMIP6 output (building on the discussion in Section 4.3).

A clear limitation of this study is the relative lack of examination of carbon cycle behavior and carbon cycle related metrics. Given the importance of the carbon cycle for emissions-driven projections, this is another clear area for future work. In the limited examination we have performed, we chose to focus on emissions-driven simulations. This choice provides the cleanest comparison between RCMs and CMIP6 models, given that many RCMs do not separate the land and ocean carbon pools, although it limits us to a relatively small set of CMIP6-comparison data (given that only few emissions-driven simulations (Jones et al., 2016) have been run by CMIP6 models). An increase in the number of emissions-driven CMIP6 ESM model output, particularly for mitigation scenarios, would greatly aid such evaluations. Using the concentration-driven simulations in future work will also provide a greater set of comparison data and will facilitate evaluation of RCMs' land and ocean carbon cycles under more varied scenarios.

Finally, given how RCMs are typically used by WG3, it appears that a truly thorough evaluation would need to consider a larger set of individual steps in the emissions-climate change cause-effect chain. Such an evaluation would provide insights into the drivers of differences between future projections based on the concentration-driven experiments typical of CMIP and results based on the all greenhouse gas emissions-driven experiments required by WG3. While it is not completely clear to us which components would need to be considered (and which could be ignored), a first suggestion of important components is: the carbon cycle, other earth system feedbacks, e.g., representation of permafrost, representation of aerosols, non-CO<sub>2</sub> greenhouse gas cycles, translation between changes in greenhouse gas concentrations and effective radiative forcing, ozone representation, land-use change albedo representation, temperature response to effective radiative forcing, and all the feedbacks and interactions. To see the full picture, a broad range of literature would need to be considered as a validation source and a wide range of experiments, spanning historical, scenario-based, and idealized experiments, would need to be performed. In performing a more thorough evaluation, an updated evaluation technique may be required. Specifically, using percentage differences from the assessed range will lead to problems when the assessed range is close to or spans zero. Hence, more sophisticated ways of evaluating the agreement between model results and assessed ranges may be required. For reasons of scope, we have not achieved such a thorough evaluation here, but we hope that this work provides a basis upon which future work can aim for the lofty goal of more complete evaluation of all of the relevant parts of the climate system.

## 6. Conclusions

We have found that the best performing RCMs can match our proxy assessment across a range of climate metrics. However, no RCM matched the proxy assessment across all metrics. At the same time, all RCMs matched the proxy assessment well for at least one metric.

Our evaluation is the first multimodel comparison of probabilistic projections from RCMs. This exercise provides a unique insight into RCMs probabilistic parameter ensembles, specifically how they compare with a set of proxy assessed ranges, which reflect wider scientific understanding of key climate metrics, and the implications of differences in probabilistic distributions for climate projections across a range of climate variables and scenarios.

Notably, although unsurprisingly, we found that models whose probabilistic distribution were constrained to the proxy assessed ranges were better able to reflect the proxy assessed ranges. This point is notable because it makes clear that if RCMs are to be used as integrators of knowledge, conveying multiple lines of evidence from one domain to another (e.g., IPCC WG1 to IPCC WG3), then RCMs whose probabilistic distributions have been constrained to the intended lines of evidence are likely to be the best tool.

Even among models which had similar levels of agreement with the proxy assessment, some divergence in future projections was observed. Given the various model structures that the reduced complexity models employ, ranging from linearized impulse response functions to 50-layer ocean models, it is not surprising that models may diverge in scenarios that go significantly beyond the domain of the validation data. Adding constraints on future performance, i.e., extending the domain of validation data (e.g., based on an independent assessment of warming in a limited subset of scenarios) would likely reduce the divergence, although such extra constraints should be carefully considered given that they risk artificially narrowing projection uncertainty.

While exercises such as the one performed here can provide helpful information about where the biases may lie, they cannot provide definitive answers about what the future holds. It is possible to make judgments about what is more reasonable based on the evaluation performed here, and to rule out clearly incorrect projections, yet it must be recognized that a definitive answer is impossible: we will not know which projections are correct until we get there, by which time it is too late for climate policy. Hence, while it is important to continue to evaluate and improve our models to remove as many sources of error as possible, it is also important that research into decision making under uncertainty (e.g., Dittrich et al., 2016; Weaver et al., 2013) continues to develop and be used because the uncertainty in projections will not disappear anytime soon, never in fact. In addition, those who use RCMs for climate projections should carefully consider how they are going to use the RCMs and how they are going to validate them before making conclusions about the implications of their projections.

In addition, we found that many of the RCMs did not reproduce the high warming seen in CMIP6 models. However, studies which constrain CMIP6 models based on observational constraints also exclude such high warming which suggests that the lack of high warming is due to the constraining applied to the RCMs, rather than structural differences between RCMs and CMIP6 models. Beyond the question of temperature projections, we found that the prescribed CO<sub>2</sub> concentrations used in the CMIP6 SSP-based experiments are at the high-end of projections made with historically constrained carbon cycles. Although, further investigations into carbon cycle behavior are required to provide a clearer picture of the influence of carbon cycle uncertainties on emissions-driven projections. Finally, we observed that a change in reference period significantly altered how well some models agreed with observations, reinforcing the need to consider more than one reference period when evaluating models.

With sufficient validation, RCMs provide a unique synthesis tool to integrate the latest scientific understanding, including its uncertainties, along the complex cause-effect chain from emissions to global-mean temperatures. Integrating this understanding in an internally consistent RCM framework, with all the implicit cross-correlations, is our best method to inform decision making and other scientific domains, for example the likelihood of exceeding a given global-mean temperature threshold under a specific emissions scenario. Further developing these tools opens vast opportunities to go beyond global-mean variables and temperature changes, and to robustly represent the complex science beneath.

### Data Availability Statement

Data and code to produce all the figures and tables can be obtained from <https://doi.org/10.5281/zenodo.4269711>.

### References

Arora, V. K., Katavouta, A., Williams, R. G., Jones, C. D., Brovkin, V., Friedlingstein, P., et al. (2020). Carbon-concentration and carbon-climate feedbacks in CMIP6 models and their comparison to CMIP5 models. *Biogeosciences*, 17(16), 4173–4222. <https://doi.org/10.5194/bg-17-4173-2020>

### Acknowledgments

Z. Nicholls benefited from support provided by the ARC Centre of Excellence for Climate Extremes (CE170100023). K. Dorheim's, D. L. Woodard's, A. Shiklomanov's, S. J. Smith's, and B. Vega-Westhoff's participation was supported by the U.S. Department of Energy, Office of Science, as part of research in MultiSector Dynamics, Earth and Environmental System Modeling Program. C. J. Smith was supported by a NERC/IIASA Collaborative Research Fellowship (NE/T009381/1). X. Su and J. Tsutsui supported by the Integrated Research Program for Advancing Climate Models (TOUGOU), Grant Number JPMXD0717935457, from the Ministry of Education, Culture, Sports, Science and Technology (MEXT), Japan. L.A. McBride's, A. P. Hope's, and R. J. Salawitch's participation was supported by the U.S. National Aeronautics and Space Administration Climate Indicators and Data Products for Future National Climate Assessments program (NX16AG34G). OSCAR development and simulations performed by T. Gasser and Y. Quilcaille were supported by the European Union's Horizon 2020 research and innovation programme under grant agreement No 820829 (CONSTRAIN project).

- Brunner, L., Pendergrass, A. G., Lehner, F., Merrifield, A. L., Lorenz, R., & Knutti, R. (2020). Reduced global warming from cmip6 projections when weighting models by performance and independence. *Earth System Dynamics Discussions*, *11*, 995–1012. <https://doi.org/10.5194/esd-11-995-2020>
- Canty, T., Mascioli, N. R., Smarte, M. D., & Salawitch, R. J. (2013). An empirical model of global climate—Part 1: A critical evaluation of volcanic cooling. *Atmospheric Chemistry and Physics*, *13*(8), 3997–4031. <https://doi.org/10.5194/acp-13-3997-2013>
- Clarke, L., Jiang, K., Akimoto, K., Babiker, M., Blanford, G., Fisher-Vanden, K., et al. (2014). Assessing transformation pathways. In O. Edenhofer, R. Pichs-Madruga, Y. Sokona, J. C. Minx, E. Farahani, S. Kadner, et al. (Eds.), *Climate change 2014: Mitigation of climate change: Contribution of working group III to the fifth assessment report of the intergovernmental panel on climate change* (pp. 413–510). Cambridge University Press.
- Dittrich, R., Wreford, A., & Moran, D. (2016). A survey of decision-making approaches for climate change adaptation: Are robust methods the way forward? *Ecological Economics*, *122*, 79–89. <https://doi.org/10.1016/j.ecolecon.2015.12.006>
- Etminan, M., Myhre, G., Highwood, E. J., & Shine, K. P. (2016). Radiative forcing of carbon dioxide, methane, and nitrous oxide: A significant revision of the methane radiative forcing. *Geophysical Research Letters*, *43*, 12614–12623. <https://doi.org/10.1002/2016GL071930>
- Eyring, V., Bony, S., Meehl, G. A., Senior, C. A., Stevens, B., Stouffer, R. J., & Taylor, K. E. (2016). *Overview of the coupled model intercomparison project phase 6 (CMIP6) experimental design and organization* (LLNL-JRNL-736881, pp. 9). Geoscientific Model Development (Online).
- Forster, P. M., Richardson, T., Maycock, A. C., Smith, C. J., Samset, B. H., Myhre, G., & Schulz, M. (2016). Recommendations for diagnosing effective radiative forcing from climate models for cmip6. *Journal of Geophysical Research: Atmospheres*, *121*, 12460–12475. <https://doi.org/10.1002/2016JD025320>
- Friedlingstein, P., Meinshausen, M., Arora, V. K., Jones, C. D., Anav, A., Liddicoat, S. K., & Knutti, R. (2014). Uncertainties in CMIP5 Climate Projections due to Carbon Cycle Feedbacks. *Journal of Climate*, *27*(2), 511–526. <https://doi.org/10.1175/JCLI-D-12-00579.1>
- Gasser, T., Ciais, P., Boucher, O., Quilcaille, Y., Tortora, M., Bopp, L., & Hauglustaine, D. (2017). The compact earth system model oscar v2.2: description and first results. *Geoscientific Model Development*, *10*(1), 271–319. <https://doi.org/10.5194/gmd-10-271-2017>
- Gasser, T., Crepin, L., Quilcaille, Y., Houghton, R. A., Ciais, P., & Obersteiner, M. (2020a). Historical CO<sub>2</sub> emissions from land use and land cover change and their uncertainty. *Biogeosciences*, *17*(15), 4075–4101. <https://doi.org/10.5194/bg-17-4075-2020>
- Gasser, T., Kechiar, M., Ciais, P., Burke, E. J., Kleinen, T., Zhu, D., et al. (2018). Path-dependent reductions in CO<sub>2</sub> emission budgets caused by permafrost carbon release. *Nature Geoscience*, *11*(11), 830–835. <https://doi.org/10.1038/s41561-018-0227-0>
- Harmsen, M. J. H. M., van Vuuren, D. P., van den Berg, M., Hof, A. F., Hope, C., Krey, V., et al. (2015). How well do integrated assessment models represent non-CO<sub>2</sub> radiative forcing? *Climatic Change*, *133*(4), 565–582. <https://doi.org/10.1007/s10584-015-1485-0>
- Haustein, K., Allen, M. R., Forster, P. M., Otto, F. E. L., Mitchell, D. M., Matthews, H. D., & Frame, D. J. (2017). A real-time Global Warming Index. *Scientific Reports*, *7*(1), 15417. <https://doi.org/10.1038/s41598-017-14828-5>
- Herger, N., Abramowitz, G., Sherwood, S., Knutti, R., Angéil, O., & Sisson, S. A. (2019). Ensemble optimisation, multiple constraints and overconfidence: A case study with future Australian precipitation change. *Climate Dynamics*, *53*(3), 1581–1596. <https://doi.org/10.1007/s00382-019-04690-8>
- Hooss, G., Voss, R., Hasselmann, K., Maier-Reimer, E., & Joos, F. (2001). A nonlinear impulse response model of the coupled carbon cycle-climate system (NICCS). *Climate Dynamics*, *18*(3–4), 189–202. <https://doi.org/10.1007/s003820100170>
- Hope, A. P., Canty, T. P., Salawitch, R. J., Tribett, W. R., & Bennett, B. F. (2017). Forecasting global warming. In *Paris climate agreement: Beacon of hope* (pp. 51–113). Springer Climate. [https://doi.org/10.1007/978-3-319-46939-3\\_2](https://doi.org/10.1007/978-3-319-46939-3_2)
- Hope, A. P., McBride, L. A., Canty, T. P., Bennett, B. F., Tribett, W. R., & Salawitch, R. J. (2020). Examining the human influence on global climate using an empirical model. *Earth and Space Science Open Archive*, *79*. <https://doi.org/10.1002/essoar.10504179.1>
- Huppmann, D., Rogelj, J., Kriegler, E., Krey, V., & Riahi, K. (2018). A new scenario resource for integrated 1.5°C research. *Nature Climate Change*, *8*(12), 1027–1030. <https://doi.org/10.1038/s41558-018-0317-4>
- Jones, C. D., Arora, V., Friedlingstein, P., Bopp, L., Brovkin, V., Dunne, J., et al. (2016). C4MIP—The Coupled Climate-Carbon Cycle Model Intercomparison Project: Experimental protocol for CMIP6. *Geoscientific Model Development*, *9*(8), 2853–2880. <https://doi.org/10.5194/gmd-9-2853-2016>
- Joos, F., Bruno, M., Fink, R., Siegenthaler, U., Stocker, T. F., Le Quéré, C., & Sarmiento, J. L. (1996). An efficient and accurate representation of complex oceanic and biospheric models of anthropogenic carbon uptake. *Tellus B: Chemical and Physical Meteorology*, *48*(3), 394–417. <https://doi.org/10.3402/tellusb.v48i3.15921>
- Kawamiya, M., Hajima, T., Tachiiri, K., Watanabe, S., & Yokohata, T. (2020). Two decades of Earth system modeling with an emphasis on Model for Interdisciplinary Research on Climate (MIROC). *Progress in Earth and Planetary Science*, *7*(1), 64. <https://doi.org/10.1186/s40645-020-00369-5>
- Leach, N. J., Jenkins, S., Nicholls, Z., Smith, C. J., Lynch, J., Cain, M., et al. (2020). Fairv2.0.0: A generalised impulse-response model for climate uncertainty and future scenario exploration. *Geoscientific Model Development Discussions*, *2020*, 1–48. <https://doi.org/10.5194/gmd-2020-390>
- Liang, Y., Gillett, N. P., & Monahan, A. H. (2020). Climate model projections of 21st century global warming constrained using the observed warming trend. *Geophysical Research Letters*, *47*, e2019GL086757. <https://doi.org/10.1029/2019GL086757>
- McBride, L. A., Hope, A. P., Canty, T. P., Bennett, B. F., Tribett, W. R., & Salawitch, R. J. (2021). Comparison of CMIP6 historical climate simulations and future projected warming to an empirical model of global climate. *Earth System Dynamics Discussions*, *12*(2), 545–579. <https://doi.org/10.5194/esd-12-545-2021>
- Meinshausen, M., Meinshausen, N., Hare, W., Raper, S. C. B., Frieler, K., Knutti, R., et al. (2009). Greenhouse-gas emission targets for limiting global warming to 2°C. *Nature*, *458*(7242), 1158–1162. <https://doi.org/10.1038/nature08017>
- Meinshausen, M., Nicholls, Z., Lewis, J., Gidden, M. J., Vogel, E., Freund, M., et al. (2020). The shared socio-economic pathway (SSP) greenhouse gas concentrations and their extensions to 2500. *Geoscientific Model Development*, *13*(8), 3571–3605. <https://doi.org/10.5194/gmd-13-3571-2020>
- Meinshausen, M., Raper, S. C. B., & Wigley, T. M. L. (2011). Emulating coupled atmosphere-ocean and carbon cycle models with a simpler model, MAGICC6 - Part 1: Model description and calibration. *Atmospheric Chemistry and Physics*, *11*(4), 1417–1456. <https://doi.org/10.5194/acp-11-1417-2011>
- Millar, R. J., Nicholls, Z., Friedlingstein, P., & Allen, M. R. (2017). A modified impulse-response representation of the global near-surface air temperature and atmospheric concentration response to carbon dioxide emissions. *Atmospheric Chemistry and Physics*, *17*(11), 7213–7228. <https://doi.org/10.5194/acp-17-7213-2017>
- Millar, R. J., Otto, A., Forster, P. M., Lowe, J. A., Ingram, W. J., & Allen, M. R. (2015). Model structure in observational constraints on transient climate response. *Climatic Change*, *131*(2), 199–211. <https://doi.org/10.1007/s10584-015-1384-4>

- Morice, C. P., Kennedy, J. J., Rayner, N. A., & Jones, P. D. (2012). Quantifying uncertainties in global and regional temperature change using an ensemble of observational estimates: The HadCRUT4 data set. *Journal of Geophysical Research*, *117*, D08101. <https://doi.org/10.1029/2011JD017187>
- Morice, C. P., Kennedy, J. J., Rayner, N. A., Winn, J. P., Hogan, E., Killick, R. E., et al. (2021). An updated assessment of near-surface temperature change from 1850: The HadCRUT5 data set. *Journal of Geophysical Research: Atmospheres*, *126*, e2019JD032361. <https://doi.org/10.1029/2019JD032361>
- Murphy, J. M., Booth, B. B. B., Boulton, C. A., Clark, R. T., Harris, G. R., Lowe, J. A., & Sexton, D. M. H. (2014). Transient climate changes in a perturbed parameter ensemble of emissions-driven earth system model simulations. *Climate Dynamics*, *43*(9), 2855–2885. <https://doi.org/10.1007/s00382-014-2097-5>
- Nicholls, Z., & Lewis, J. (2021). *Reduced Complexity Model Intercomparison Project (RCMIP) protocol*. Zenodo. <https://doi.org/10.5281/zenodo.4589756>
- Nicholls, Z., Meinshausen, M., Lewis, J., Gieseke, R., Dommenges, D., Dorheim, K., et al. (2020). Reduced complexity model intercomparison project phase 1: Introduction and evaluation of global-mean temperature response. *Geoscientific Model Development*, *13*(11), 5175–5190. <https://doi.org/10.5194/gmd-13-5175-2020>
- O'Neill, B. C., Tebaldi, C., van Vuuren, D. P., Eyring, V., Friedlingstein, P., Hurtt, G., et al. (2016). The scenario model intercomparison project (ScenarioMIP) for CMIP6. *Geoscientific Model Development*, *9*(9), 3461–3482. <https://doi.org/10.5194/gmd-9-3461-2016>
- Rogelj, J., Meinshausen, M., & Knutti, R. (2012). Global warming under old and new scenarios using IPCC climate sensitivity range estimates. *Nature Climate Change*, *2*(4), 248–253. <https://doi.org/10.1038/nclimate1385>
- Rogelj, J., Shindell, D., Jiang, K., Fifita, S., Forster, P., Ginzburg, V., et al. (2018). Mitigation pathways compatible with 1.5°C in the context of sustainable development. In G. Flato, J. Fuglestedt, R. Mrabet, & R. Schaeffer (Eds.), *Global warming of 1.5°C: An IPCC special report on the impacts of global warming of 1.5°C above pre-industrial levels and related global greenhouse gas emission pathways, in the context of strengthening the global response to the threat of climate change, sustainable development, and efforts to eradicate poverty* (pp. 93–174). IPCC/WMO. Retrieved from <http://www.ipcc.ch/report/sr15/>
- Sanderson, B. M., Wehner, M., & Knutti, R. (2017). Skill and independence weighting for multi-model assessments. *Geoscientific Model Development*, *10*(6), 2379–2395. <https://doi.org/10.5194/gmd-10-2379-2017>
- Schlesinger, M. E., Jiang, X., & Charlson, R. J. (1992). Implication of anthropogenic atmospheric sulfate for the sensitivity of the climate system. In L. Rosen & R. Glasser (Eds.), *Climate change and energy policy: Proceedings of the international conference on global climate change: Its mitigation through improved production and use of energy*. (pp. 75–108). New York: American Institute of Physics
- Schwarber, A. K., Smith, S. J., Hartin, C. A., Vega-Westhoff, B. A., & Srivier, R. (2019). Evaluating climate emulation: Fundamental impulse testing of simple climate models. *Earth System Dynamics*, *10*(4), 729–739. <https://doi.org/10.5194/esd-10-729-2019>
- Sherwood, S. C., Webb, M. J., Annan, J. D., Armour, K. C., Forster, P. M., Hargreaves, J. C., et al. (2020). An assessment of Earth's climate sensitivity using multiple lines of evidence. *Reviews of Geophysics*, *58*, e2019RG000678. <https://doi.org/10.1029/2019RG000678>
- Simmons, A. J., Berrisford, P., Dee, D. P., Hersbach, H., Hirahara, S., & Thépaut, J.-N. (2017). A reassessment of temperature variations and trends from global reanalyses and monthly surface climatological datasets. *Quarterly Journal of the Royal Meteorological Society*, *143*(702), 101–119. <https://doi.org/10.1002/qj.2949>
- Skeie, R. B., Bernsten, T., Aldrin, M., Holden, M., & Myhre, G. (2018). Climate sensitivity estimates—Sensitivity to radiative forcing time series and observational data. *Earth System Dynamics*, *9*(2), 879–894. <https://doi.org/10.5194/esd-9-879-2018>
- Skeie, R. B., Fuglestedt, J., Bernsten, T., Peters, G. P., Andrew, R., Allen, M., & Kallbekken, S. (2017). Perspective has a strong effect on the calculation of historical contributions to global warming. *Environmental Research Letters*, *12*(2), 024022. <https://doi.org/10.1088/1748-9326/aa5b0a>
- Skeie, R. B., Peters, G. P., Fuglestedt, J., & Andrew, R. (2021). A future perspective of historical contributions to climate change. *Climatic Change*, *164*(1), 24. <https://doi.org/10.1007/s10584-021-02982-9>
- Smith, C. J., Forster, P. M., Allen, M., Leach, N., Millar, R. J., Passerello, G. A., & Regayre, L. A. (2018). FAIR v1.3: A simple emissions-based impulse response and carbon cycle model. *Geoscientific Model Development*, *11*(6), 2273–2297. <https://doi.org/10.5194/gmd-11-2273-2018>
- Smith, C. J., Kramer, R. J., Myhre, G., Alterskjær, K., Collins, W., Sima, A., et al. (2020). Effective radiative forcing and adjustments in cmip6 models. *Atmospheric Chemistry and Physics*, *20*(16), 9591–9618. <https://doi.org/10.5194/acp-20-9591-2020>
- Smith, C. J., Kramer, R. J., Myhre, G., Forster, P. M., Soden, B. J., Andrews, T., & Watson-Parris, D. (2018). Understanding rapid adjustments to diverse forcing agents. *Geophysical Research Letters*, *45*, 12023–12031. <https://doi.org/10.1029/2018GL079826>
- Su, X., Shioyama, H., Tanaka, K., Fujimori, S., Hasegawa, T., Hijioka, Y., et al. (2018). How do climate-related uncertainties influence 2 and 1.5°C pathways? *Sustainability Science*, *13*(2), 291–299. <https://doi.org/10.1007/s11625-017-0525-2>
- Su, X., Tachiiri, K., Tanaka, K., Watanabe, M., & Kawamiya, M. (2020). *Source attributions of radiative forcing by regions, sectors, and climate forcings*. arXiv preprint. arXiv:2009.07472.
- Su, X., Takahashi, K., Fujimori, S., Hasegawa, T., Tanaka, K., Kato, E., et al. (2017). Emission pathways to achieve 2.0°C and 1.5°C climate targets. *Earth's Future*, *5*, 592–604. <https://doi.org/10.1002/2016EF000492>
- Tokarska, K. B., Stolpe, M. B., Sippel, S., Fischer, E. M., Smith, C. J., Lehner, F., & Knutti, R. (2020). Past warming trend constrains future warming in cmip6 models. *Science Advances*, *6*(12), eaaz9549. <https://doi.org/10.1126/sciadv.aaz9549>
- Tsutsui, J. (2017). Quantification of temperature response to CO2 forcing in atmosphere-ocean general circulation models. *Climatic Change*, *140*(2), 287–305. <https://doi.org/10.1007/s10584-016-1832-9>
- Tsutsui, J. (2020). Diagnosing transient response to CO2 forcing in coupled atmosphere-ocean model experiments using a climate model emulator. *Geophysical Research Letters*, *47*, e2019GL085844. <https://doi.org/10.1029/2019GL085844>
- Uhe, P., Otto, F. E., Rashid, M. M., & Wallom, D. C. (2016). Utilizing amazon web services to provide an on demand urgent computing facility for climateprediction.net. In *2016 IEEE 12th international conference on e-science (e-science)* (pp. 407–413). IEEE.
- van Vuuren, D. P., Lowe, J., Stehfest, E., Gohar, L., Hof, A. F., Hope, C., et al. (2011). How well do integrated assessment models simulate climate change? *Climatic Change*, *104*(2), 255–285. <https://doi.org/10.1007/s10584-009-9764-2>
- Vega-Westhoff, B., Srivier, R. L., Hartin, C. A., Wong, T. E., & Keller, K. (2019). Impacts of observational constraints related to sea level on estimates of climate sensitivity. *Earth's Future*, *7*(6), 677–690. <https://doi.org/10.1029/2018ef001082>
- von Schuckmann, K., Cheng, L., Palmer, M. D., Hansen, J., Tassone, C., Aich, V., et al. (2020). Heat stored in the earth system: Where does the energy go? *Earth System Science Data*, *12*(3), 2013–2041. <https://doi.org/10.5194/essd-12-2013-2020>
- Weaver, C. P., Lempert, R. J., Brown, C., Hall, J. A., Revell, D., & Sarewitz, D. (2013). Improving the contribution of climate model information to decision making: The value and demands of robust decision frameworks. *WIREs Climate Change*, *4*(1), 39–60. <https://doi.org/10.1002/wcc.202>

## References From the Supporting Information

- Arora, V. K., Boer, G. J., Friedlingstein, P., Eby, M., Jones, C. D., Christian, J. R., et al. (2013). Carbon-concentration and carbon-climate feedbacks in CMIP5 Earth system models. *Journal of Climate*, *26*(15), 5289–5314. <https://doi.org/10.1175/JCLI-D-12-00494.1>
- Bellouin, N., Davies, W., Shine, K. P., Quaas, J., Miltenstädt, J., Forster, P. M., et al. (2020). Radiative forcing of climate change from the copernicus reanalysis of atmospheric composition. *Earth System Science Data*, *12*(3), 1649–1677. [10.5194/essd-12-1649-2020](https://doi.org/10.5194/essd-12-1649-2020)
- Cowan, K., & Way, R. G. (2014). Coverage bias in the hadcrut4 temperature series and its impact on recent temperature trends. *Quarterly Journal of the Royal Meteorological Society*, *140*(683), 1935–1944. <https://doi.org/10.1002/qj.2297>
- Cummins, D. P., Stephenson, D. B., & Stott, P. A. (2020). Optimal estimation of stochastic energy balance model parameters. *Journal of Climate*, *33*(18), 7909–7926.
- Gasser, T., Crepin, L., Quilcaille, Y., Houghton, R. A., Ciais, P., & Obersteiner, M. (2020). Historical CO<sub>2</sub> emissions from land use and land cover change and their uncertainty. *Biogeosciences*, *17*(15), 4075–4101. <https://doi.org/10.5194/bg-17-4075-2020>
- Gasser, T., Peters, G. P., Fuglestedt, J. S., Collins, W. J., Shindell, D. T., & Ciais, P. (2017). Accounting for the climate-carbon feedback in emission metrics. *Earth System Dynamics*, *8*(2), 235–253. <https://doi.org/10.5194/esd-8-235-2017>
- Geoffroy, O., Saint-Martin, D., Bellon, G., Voldoire, A., Olivie, D. J. L., & Tyteca, S. (2013a). Transient climate response in a two-layer energy-balance model. Part II: Representation of the efficacy of deep-ocean heat uptake and validation for CMIP5 AOGCMs. *Journal of Climate*, *26*(6), 1859–1876. <https://doi.org/10.1175/JCLI-D-12-00196.1>
- Geoffroy, O., Saint-Martin, D., Olivie, D. J. L., Voldoire, A., Bellon, G., & Tyteca, S. (2013b). Transient climate response in a two-layer energy-balance model. part i: Analytical solution and parameter calibration using CMIP5 AOGCM experiments. *Journal of Climate*, *26*(6), 1841–1857. <https://doi.org/10.1175/jcli-d-12-00195.1>
- Hartin, C. A., Patel, P., Schwarber, A., Link, R. P., & Bond-Lamberty, B. P. (2015). A simple object-oriented and open-source model for scientific and policy analyses of the global climate system—hector v1.0. *Geoscientific Model Development*, *8*(4), 939–955. <https://doi.org/10.5194/gmd-8-939-2015>
- Hauglustaine, D. A., Balkanski, Y., & Schulz, M. (2014). A global model simulation of present and future nitrate aerosols and their direct radiative forcing of climate. *Atmospheric Chemistry and Physics*, *14*(20), 11031–11063. <https://doi.org/10.5194/acp-14-11031-2014>
- IPCC (2013). Summary for policymakers [Book Section]. In T. Stocker 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* (pp. 1–30). Cambridge, UK: Cambridge University Press. <https://doi.org/10.1017/CBO9781107415324.004>
- Kriegler, E. (2005). Imprecise probability analysis for integrated assessment of climate change (Doctoral thesis, Universität Potsdam). Retrieved from <https://publishup.uni-potsdam.de/opus4-ubp/frontdoor/index/index/docId/497>
- Myhre, G., Shindell, D., Bréon, F.-M., Collins, W., Fuglestedt, J., Huang, J., et al. (2013). Anthropogenic and natural radiative forcing [Book Section]. In T. Stocker 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* (pp. 659–740). Cambridge, UK: Cambridge University Press. <https://doi.org/10.1017/CBO9781107415324.018>
- Tanaka, K., Kriegler, E., Bruckner, T., Hooss, G., Knorr, W., Raddatz, T., & Tol, R. (2007). Aggregated carbon cycle, atmospheric chemistry and climate model (acc2): Description of forward and inverse mode. Retrieved from [https://pure.mpg.de/rest/items/item\\_994422/component/file\\_994421/content](https://pure.mpg.de/rest/items/item_994422/component/file_994421/content)
- Tsutsui, J. (2021). Minimal CMIP Emulator (MCE v1.2). Zenodo. <https://doi.org/10.5281/zenodo.4604695>
- Urban, N. M., Holden, P. B., Edwards, N. R., Sriver, R. L., & Keller, K. (2014). Historical and future learning about climate sensitivity. *Geophysical Research Letters*, *41*, 2543–2552. <https://doi.org/10.1002/2014GL059484>