



Detection and Attribution of Climate Change Using a Neural Network

Constantin Bône, Guillaume Gastineau, Sylvie Thiria, Patrick Gallinari,
Carlos Mejia

► To cite this version:

Constantin Bône, Guillaume Gastineau, Sylvie Thiria, Patrick Gallinari, Carlos Mejia. Detection and Attribution of Climate Change Using a Neural Network. Journal of Advances in Modeling Earth Systems, 2023, 15 (10), 10.1029/2022ms003475 . hal-04355503

HAL Id: hal-04355503

<https://hal.science/hal-04355503>

Submitted on 21 Dec 2023

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution - NonCommercial - NoDerivatives 4.0
International License



RESEARCH ARTICLE

10.1029/2022MS003475

Key Points:

- We present a non linear method based on neural network to attribute the global mean surface air temperature variability to different forcings
- We use a convolutional neural network associated with a backward optimization to estimate the climate response to the different external forcings
- The attributable forcings are consistent with those obtained using another state-of-the-art method

Supporting Information:

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

Correspondence to:

C. Bône,
constantin.bone@sorbonne-universite.fr

Citation:

Bône, C., Gastineau, G., Thiria, S., Gallinari, P., & Mejia, C. (2023). Detection and attribution of climate change using a neural network. *Journal of Advances in Modeling Earth Systems*, 15, e2022MS003475. <https://doi.org/10.1029/2022MS003475>

Received 25 OCT 2022

Accepted 23 AUG 2023

Author Contributions:

Conceptualization: Constantin Bône, Guillaume Gastineau, Sylvie Thiria, Patrick Gallinari, Carlos Mejia

Data curation: Guillaume Gastineau

Methodology: Constantin Bône, Guillaume Gastineau, Sylvie Thiria, Patrick Gallinari, Carlos Mejia

Resources: Guillaume Gastineau

Software: Constantin Bône, Carlos Mejia

Supervision: Guillaume Gastineau, Sylvie Thiria

© 2023 The Authors. Journal of Advances in Modeling Earth Systems 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.

Detection and Attribution of Climate Change Using a Neural Network

Constantin Bône^{1,2} , Guillaume Gastineau¹, Sylvie Thiria¹, Patrick Gallinari^{2,3}, and Carlos Mejia¹

¹UMR LOCEAN, IPSL, Sorbonne Université, IRD, CNRS, MNHN, Paris, France, ²UMR ISIR, Sorbonne Université, CNRS, INSERM, Paris, France, ³Criteo AI Lab, Paris, France

Abstract A new detection and attribution method is presented and applied to the global mean surface air temperature (GSAT) from 1900 to 2014. The method aims at attributing the climate changes to the variations of greenhouse gases, anthropogenic aerosols, and natural forcings. A convolutional neural network (CNN) is trained using the simulated GSAT from historical and single-forcing simulations of 12 climate models. Then, we perform a backward optimization with the CNN to estimate the attributable GSAT changes. Such a method does not assume additivity in the effects of the forcings. The uncertainty in the attributable GSAT is estimated by sampling different starting points from single-forcing simulations and repeating the backward optimization. To evaluate this new method, the attributable GSAT changes are also calculated using the regularized optimal fingerprinting (ROF) method. Using synthetic non-additive data, we first find that the neural network-based method estimates attributable changes better than ROF. When using GSAT data from climate model, the attributable anomalies are similar for both methods, which might reflect that the influence of forcing is mainly additive for the GSAT. However, we found that the uncertainties given both methods are different. The new method presented here can be adapted and extended in future work, to investigate the non-additive changes found at the local scale or on other physical variables.

Plain Language Summary In order to design effective adaptation policies, it is essential to have reliable estimates of the effect of anthropogenic activities on the climate. For that purpose, a new attribution method based on a neural network is designed and evaluated. The method estimates the past global mean surface air temperatures anomalies caused by the changes in the greenhouse gases concentration, the variation of anthropogenic aerosols, and the variations driven by naturally occurring phenomena. To build this estimation, the data from observations and climate models are used. This methodology is compared with another state-of-the-art method. The results of both methods are evaluated and discussed. The proposed method provides better estimations in the case of large non-additivity of the causes of climate change and can be applied to other physical variables or at the regional scale. In the case of the global mean surface air temperature, the method presented provides estimation similar to other methods.

1. Introduction

Detection and attribution of climate change is key to understanding past climate change and devising adaptation policies. This problem is an important part of IPCC reports (Eyring et al., 2021) as it directly inquires about the impact of anthropogenic activities on the climate system. Detection aims to compare climate change with internal variability. A change is detected if it exceeds the anomalies generated by the internal climate variability. Internal variability refers to climate variations resulting from processes intrinsic to the climate system, occurring in the absence of external forcing. Internal variability may arise from processes within each of the climate system components (atmosphere, ocean, land surface, cryosphere) or may emerge from their interactions (Cassou et al., 2018). For instance, the global mean surface air temperature (GSAT) varies by a few tenths of degrees during the El Niño or La Niña phases of the El Niño Southern Oscillation (Neelin et al., 1998). Similarly, the Pacific decadal variability and the Atlantic multi-decadal variability can also influence the GSAT (Z. Li et al., 2020; Meehl et al., 2016). Forcing agents external to the climate system, known as external forcings, can also cause climate changes. The dominant forcings in the historical period (i.e., 1850 to present-day) are the increase in the concentration of greenhouse gases, the variations of the aerosol concentrations, the variations of incoming solar radiation, the changes in land use and stratospheric ozone concentration (Masson-Delmotte et al., 2021). Attribution then aims to explain and quantify the impacts of the different forcings. Anthropogenically driven and naturally occurring forcings are often considered separately to understand the impact of

Validation: Guillaume Gastineau, Sylvie Thiria

Writing – original draft: Constantin Bône

Writing – review & editing: Guillaume Gastineau, Sylvie Thiria, Patrick Gallinari, Carlos Mejia

human activities. Natural forcings include the effects of natural sources of aerosols and solar activity. The anthropogenic effects include the contributions of other effects. Hasselmann (1993) defined a method called “optimal fingerprinting” for detection and attribution relying on climate model simulations and observations. This method has been improved to build more reliable uncertainties and to check for the consistency between models and observations (Allen & Tett, 1999), or to account for the residual internal variability in ensembles of climate model simulations (Allen & Stott, 2003). To better account for the uncertainty in the estimation of forcings Ribes et al. (2013) proposed to use a regularized estimator of the covariance matrix of internal variability. A review, based, among other, on regularized optimal forcing estimates, concluded that the likely range (5%–95% range) of the attributable anthropogenic GSAT anomaly in 2010–2019 relative to 1850–1900 is between +0.8°C and +1.3°C (Eyring et al., 2021). The anomaly attributable to greenhouse gases reported is +1.0°C to +2.0°C, while it is from −0.8°C to 0.0°C for other anthropogenic forcings, and from −0.1°C to +0.1°C for natural forcings.

However, the optimal fingerprinting has several limitations such as the loss of information due to the reduction of the temporal and spatial dimensionality of data, needed to make a proper approximation of the covariance matrix of internal variability. Another problem is the additivity assumption where the individual forcing effects are summed together to estimate the climate response to the sum of forcings even if it is verified for the attribution of historical GSAT (Marvel et al., 2015; Shiogama et al., 2013). This additivity assumption also found to be invalid for precipitation (Marvel et al., 2015), the surface air temperature changes driven by greenhouse gases and aerosols can be non-additive over the extra-tropical regions such as the Arctic (Deng et al., 2020) or the Southern Hemisphere (Pope et al., 2020).

To take account of non-additive changes, we present here a new method for attributing past climate using machine learning. A neural network is a machine learning method consisting of consecutive hidden layers of nonlinear transformations and adjustable weights and biases which are determined by applying gradient descent using backpropagation (Goodfellow et al., 2016). It is a statistical tool increasingly used in recent years in many scientific fields (Choudhary et al., 2022). Convolutional neural networks (CNNs; Yamashita et al., 2018) are a class of non-linear neural networks used notably in imagery problems (O'Shea & Nash, 2015). Their main characteristic is the use of a learnable kernel that slides along the input data. The CNNs have also shown their great capacity to analyze time series and other one-dimensional patterns (Kiranyaz et al., 2021) and have become common machine learning tools. For instance, without being exhaustive, neural networks have been used in climate science to predict the evolution of El Nino Southern Oscillation (Ham et al., 2019), to identify storm structures (Gagne et al., 2019), for weather prediction (Gagne et al., 2019; Lam et al., 2022), or for detection studies (Barnes et al., 2019; Labe & Barnes, 2021). However, they are still emerging in large parts of the geosciences.

Here, we propose an alternative attribution framework based on a CNN to account for interactions between the forcings. To the best of our knowledge, this is the first attempt to apply a neural network to the problem of detection and attribution of climate change. We compare the results obtained with the neural-network based attribution method with those resulting from regularized optimal fingerprinting (ROF). We chose to study the GSAT as it is widely studied in the detection and attribution literature in order to properly introduce our methodology. We investigate the effects of greenhouse gases, anthropogenic aerosols and natural forcings. In the future, this attribution method based on a neural network could be applied to other physical variables such as precipitation, or changes at the regional scale where non-additivity are expected to be more important (Good et al., 2015).

To evaluate our neural network-based attribution method and compare it to ROF, we first build synthetic data to assess the ability of methods to take non-additivities into account. Then we use a perfect model approach. This consists of removing data coming from one climate model and treating its simulations as pseudo-observations. The estimated effect of each forcing is then compared to their actual simulated effects.

The article is organized as follows. In Section 2, we present the data and the preprocessing applied and how we built up synthetic data. In Section 3, we present the neural network and its direct performance. We also introduce the two attribution methods used in this paper: backward optimization and ROF. In Section 4, we present the results obtained by the two attribution methods. Finally, in Section 5, we conclude and discuss the limitations as well as future perspectives.

Table 1
Presentation of the Climate Models Used

Model	n_{GHG}	n_{AER}	n_{NAT}	n_{HIST}	PI (year)	σ_{PI} (°C)	Reference
CanESM5	50	30	30	65	1,000	0.10	Swart et al. (2019)
CESM2	3	3	2	11	500	0.13	Danabasoglu et al. (2020)
IPSL-CM6-LR	10	10	10	32	1,000	0.15	Boucher et al. (2020)
ACCESS-ESM1-5	3	3	3	30	500	0.11	Ziehn et al. (2020)
BCC-CSM2-MR	3	3	3	3	600	0.17	Wu et al. (2019)
CNRM-CM6-1	9	10	10	30	500	0.13	Voldoire et al. (2019)
FGOALS-g3	3	3	3	6	700	0.10	L. Li et al. (2020)
HadGEM3	4	4	4	5	500	0.11	Roberts et al. (2019)
MIROC6	3	3	3	50	500	0.13	Tatebe et al. (2019)
MRI-ESM2.0	5	5	5	7	500	0.10	Yukimoto et al. (2019)
NorESM2-LM	3	3	3	3	500	0.15	Seland et al. (2020)
GISS-E2-1-G	5	7	15	19	500	0.15	Kelley et al. (2020)

Note. n_{GHG} , n_{AER} , n_{NAT} , and n_{HIST} denote the number of members used for GHG, AER, NAT, and HIST. The duration of the PI simulation is indicated, in year. σ_{PI} denotes the year to year standard deviation of the GSAT from PI, in °C.

2. Model and Data

2.1. Climate Models Simulations

In this section, we present the climate model data used in this study. We use the monthly surface air temperature from the outputs of the Coupled Model Intercomparison Project 6 phase (CMIP6; Eyring et al., 2016) and of the Detection and Attribution Model Intercomparison Project (DAMIP; Gillett et al., 2016) panel of CMIP6. All simulations from CMIP6 use the same experimental protocol with identical boundary conditions based on reconstructions and observations.

We use the historical simulations, called HIST, to obtain estimation of the combined effect of the forcings. These simulations use as variable boundary conditions all external forcings from 1850 to 2014. This includes the reconstructed concentrations of greenhouse gases, anthropogenic aerosols and ozone, and the estimated past variations of solar incoming radiation and land-use.

We also use single-forcing simulations to obtain estimation of the individual effect of the forcings. These simulations use as variable boundary conditions only one of the external forcings, all the other external forcings being fixed at their value from 1850. We use the single-forcing simulations hist-GHG denoted later GHG, hist-aer denoted AER, and hist-nat denoted NAT dedicated respectively to greenhouse gas concentrations, anthropogenic aerosols, and natural forcings (i.e., volcanic aerosol and solar variations) as variable forcings for the same period (1850–2014). The effect of stratospheric ozone and land use was not investigated as only a few simulations have been performed in CMIP6, and because their effective radiative forcings are much smaller than the ones of greenhouse gases, aerosols or natural forcings (C. J. Smith et al., 2020).

We also use the preindustrial control simulations, called PI, to estimate of the effects of internal variability. These control simulations use fixed forcings from their estimated pre-industrial levels corresponding that of 1850. The PI simulations are multi-centennal with usually a single realization for each climate model. These simulations show a small drift due to incomplete spin-up or nonclosure of the energy budget (Hobbs et al., 2016). Hereafter such small long-term drift (Irving et al., 2021) is deleted from each PI simulations by removing a quadratic trend (Gupta et al., 2013) of the simulated GSAT before analysis in all simulations.

All simulations but PI includes multiple realizations called ensemble members and denoted later as members. The members use different initial conditions, which are sampled from the PI simulation. We use 12 atmosphere-ocean general circulation models (AOGCMs, see Table 1 for details) where at least two members are available for the simulations HIST, GHG, AER and NAT.

2.2. Observations

We use observations of the 2 m air temperature from HadCRUT5 (Morice et al., 2021). The gridded data is a blend of the CRUTEM5 (Osborn et al., 2021) land-surface air temperature data set and the HadSST4 (Kennedy et al., 2019) sea-surface temperature (SST) data set. Such a blending is necessary because there are few observations of temperature at 2 m over the oceans compared to SST observations. The resulting globally averaged quantity is called global mean surface temperature (GMST) and it differs from the GSAT which is solely based on surface air temperature. To correct this we multiply by 1.06 the GMST from observation to estimate the observed GSAT, as estimated by Richardson et al. (2018).

2.3. Pre-Processing

All monthly climate model data are aggregated to an annual mean and spatially averaged from 90°S to 90°N to provide the GSAT. We then estimate the temperature anomalies compared to the pre-industrial period.

We remove the time mean GSAT of PI from the GHG, AER, and NAT simulations. For observations and HIST, we compute the average temperature during the 1850–1900 period and remove it from the GSAT. Hereafter we only use the data from 1900 to 2014 period (115 years).

The simulated and observed GSAT can be separated into a forced component and an internally-generated climate variability component. To reduce the effects of internal climate variability we apply a low-pass filter to the GSAT of the GHG and AER simulations. We use a Lanczos low-pass filter (Burger & Burge, 2009), with a window size of 21 years, and a cutoff period of 10 years. The endpoints are estimated by extending the time series by replicating the mean value of the first and last 10 years of each simulation. This should not alter the estimated effect of greenhouse gases or aerosols on the GSAT as both forcings only show multi-decadal and longer fluctuations in terms of effective radiative forcing (Gulev et al., 2021). We do not apply this procedure to NAT and HIST because the emission of aerosol from volcanic eruptions induces an intense cooling for the next 2–5 years, and such smoothing would degrade the forced anomalies. This smoothing procedure only lead to minor improvements regarding the estimated uncertainties (not shown).

We illustrate in Figure 1 the processed data for all climate models, observations and the multi-model mean (MMM) for each forcing. To compute the MMM we first compute the ensemble mean (i.e., averaging all ensemble member) for each climate model and then we average the 12 ensemble means. In all models, GHG shows a monotonic warming with an increasing slope since the 1960s, as expected from the greenhouse gases emissions. In AER, the aerosol induces a cooling with a pronounced slope from the 1940s to 1980s, and a plateau from 1980 to 2014. NAT shows small cooling from 0.1°C to 0.4°C only occurring after the major eruptive volcanic eruptions of Agung (1963), El Chichon (1982) and Pinatubo (1991). HIST shows a monotonic warming less pronounced than GHG with also a cooling a few years after the major volcanic eruptions. In all simulations, the internal variability is important, as illustrated by the fluctuations visible in each member (thin lines) and is reduced in the ensemble mean (thick lines).

2.4. Synthetic Data

To investigate the performance of the attribution methods when considering external forcings with non-additive influences, a synthetic data set is generated. We generate three time series of size 115 denoted f_1 , f_2 , and f_3 , that represents the forced effects of three synthetic forcings. These time series are constructed to have similarities with the expected influence of the greenhouse gases, aerosols, and natural forcing for f_1 , f_2 , and f_3 , respectively (see Figure 2, red, green, and blue lines). However, the expressions of f_1 , f_2 , and f_3 remain arbitrary and are not meant to represent simulated or observed climate. We detail in Text S1 in Supporting Information S1 the analytic expressions used to build the time series. We construct the total effect of the three forcings combined, noted r , using two additional terms compared to the additive case:

$$r = f_1 + 0.3f_1^2 + f_2 + f_3 + 0.1f_1f_2 \quad (1)$$

Using an analogy with climate, anomalies are considered to result from the addition of a forced and an internally-generated variability component (see Figure 1). We add an additional variability to f_1 , f_2 , and f_3 and r that only represent the forced component. To generate this variability, we fit a first order autoregressive (AR1)

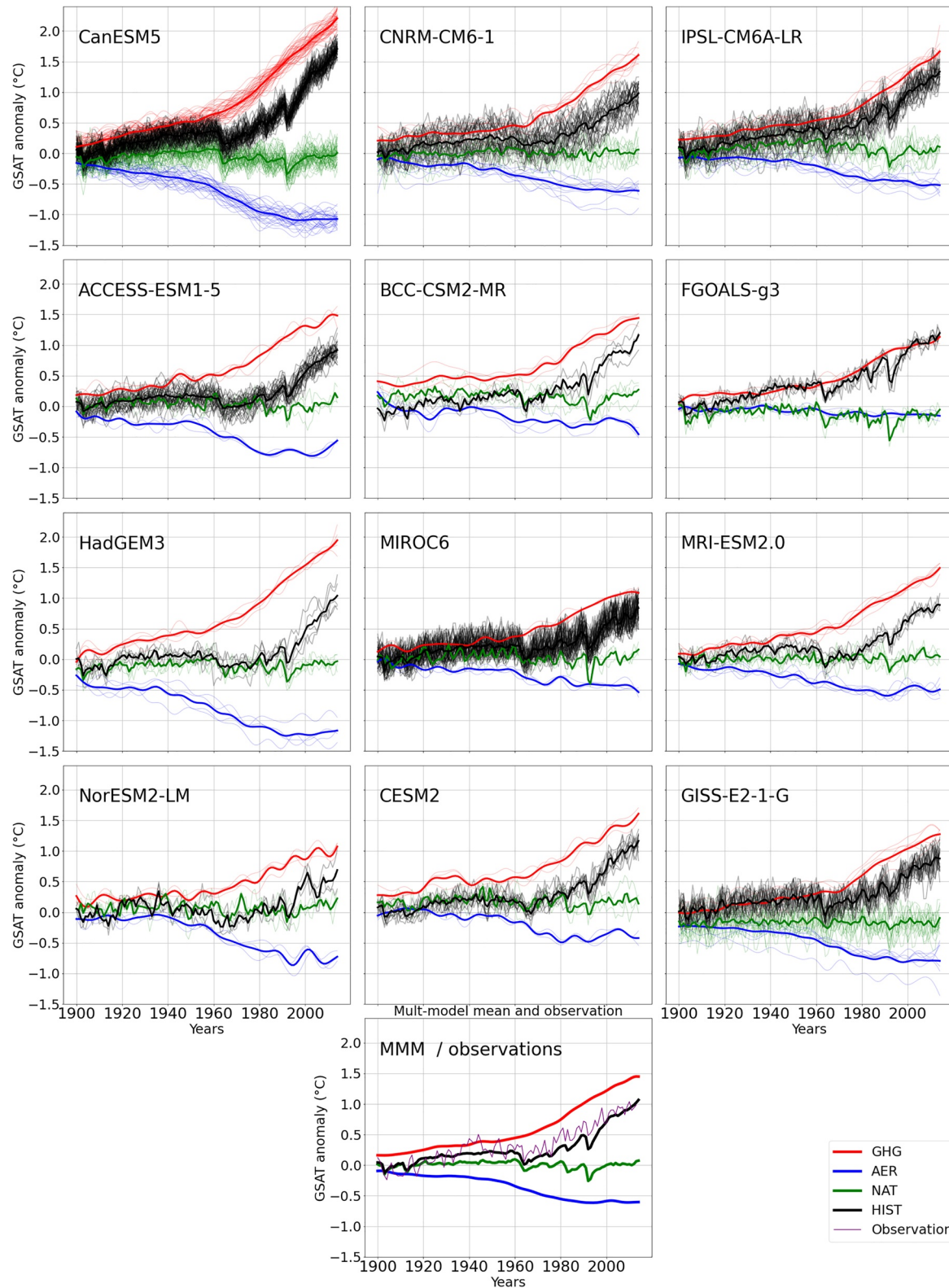


Figure 1. Global mean surface air temperature (GSAT) anomaly simulated by each model and (lower panel only) multi-model mean (MMM) and observed GSAT. Black lines show the HIST members. Red lines show the GHG members. Green lines show the NAT members. Blue lines show the AER members. The purple line shows the observations in the lower panel. Bold lines of the same colors show the ensemble mean.

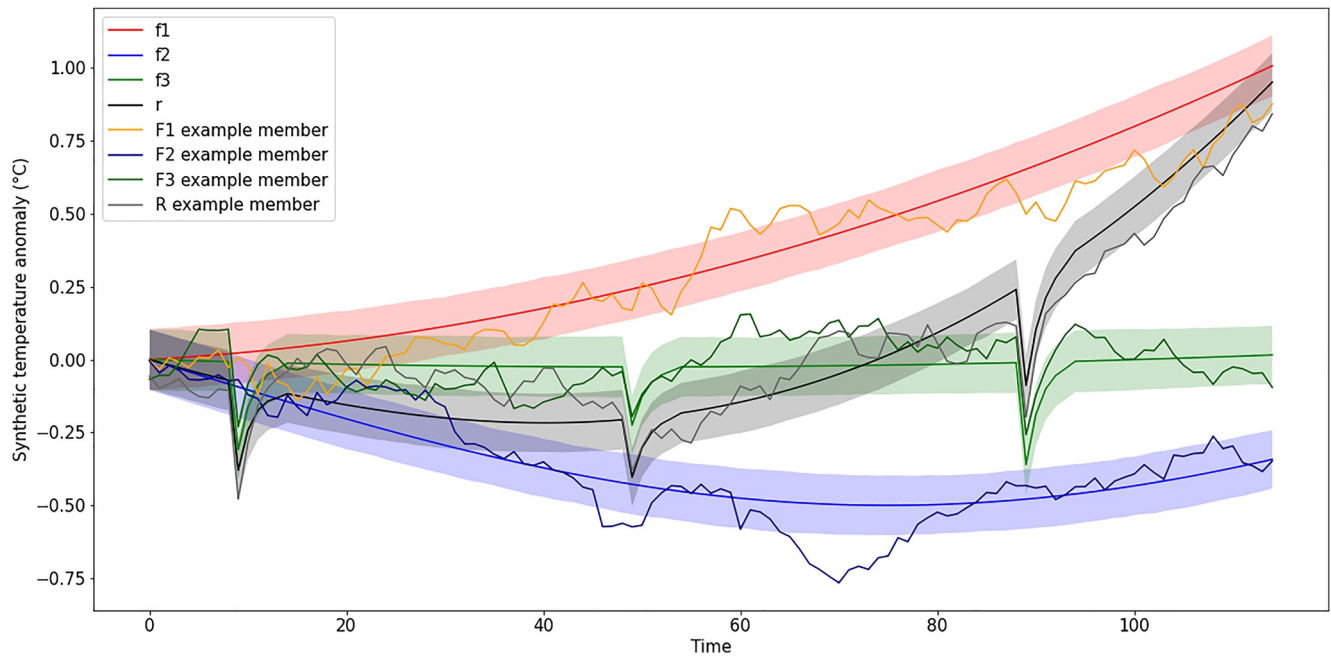


Figure 2. Synthetic time series f_1 (red), f_2 (blue), f_3 (green), and r (black). A randomly chosen time series after adding the variability is illustrated for F_1 (orange), F_2 (dark blue), F_3 (dark green) and R (gray). Colors shades indicate one standard deviation across the 100 surrogate time series obtained for each pseudo-forcings and their response.

model using the time series obtained from the concatenated PI simulations from all models. This AR1 model is then used to generate 410 surrogate time series that are added to f_1 , f_2 , f_3 , and r . This provides the 100 time series for each forcings denoted F_1 , F_2 , F_3 , and 110 time series R resulting from the combined forcings (see Figure 2).

3. Methods

3.1. Backward Optimization of a Neural Network

3.1.1. Neural Network

In this section, we describe the neural network used. We determine the relationship linking the GSAT from HIST to that of GHG, AER, and NAT using a CNN. In the training procedure, we use the GSAT from AER, GHG, and NAT as inputs and the GSAT from HIST as the target. Our goal is to construct a predictor that captures the role of all forcings combined. We assume that stratospheric ozone and land use do not affect this relationship.

A schematic of the CNN used is shown in Figure 3. CNNs can be used to construct relatively simple neural networks as the number of weights and biases is directly decided by the size and number of the filters used. We assume that this architecture is suitable in the present case the size of the data set is relatively small compared to other neural network applications. This might limit the overfitting which occurs when a neural network model performs significantly better for training data than it does for new data. In our case, a one-dimensional kernel is applied to the temporal dimension. To fix the values of the weights and biases of the convolutional layers, a neural network needs a learning data set composed of input-output pairs. The outputs are the GSAT of one HIST member while the inputs are built with one member for each single-forcing simulations. We build this data set by going through all combinations of GHG, AER, NAT, and HIST members of the same climate model. To test the backward optimization (see Section 3.3), we removed one HIST member from each climate model and 10 for the IPSL-CM6-LR model from these combinations to serve as test data set. This provides for the training of the neural network $N_d = (n_{\text{HIST}} - 1) n_{\text{GHG}} n_{\text{AER}} n_{\text{NAT}}$ 4-tuples for each climate model except for IPSL-CM6-LR with $(n_{\text{HIST}} - 10) n_{\text{GHG}} n_{\text{AER}} n_{\text{NAT}}$ 4-tuples. We note N_d the total number of the 4-tuples obtained for all models. The training data set is thus of size N_d which is of the order of 10^5 while an individual input is of size (3,115) and its corresponding output of size (1,115). The usual practice is to go through this database several times to train the CNN. However, we have altered the procedure to provide a similar weight to all models during the training.

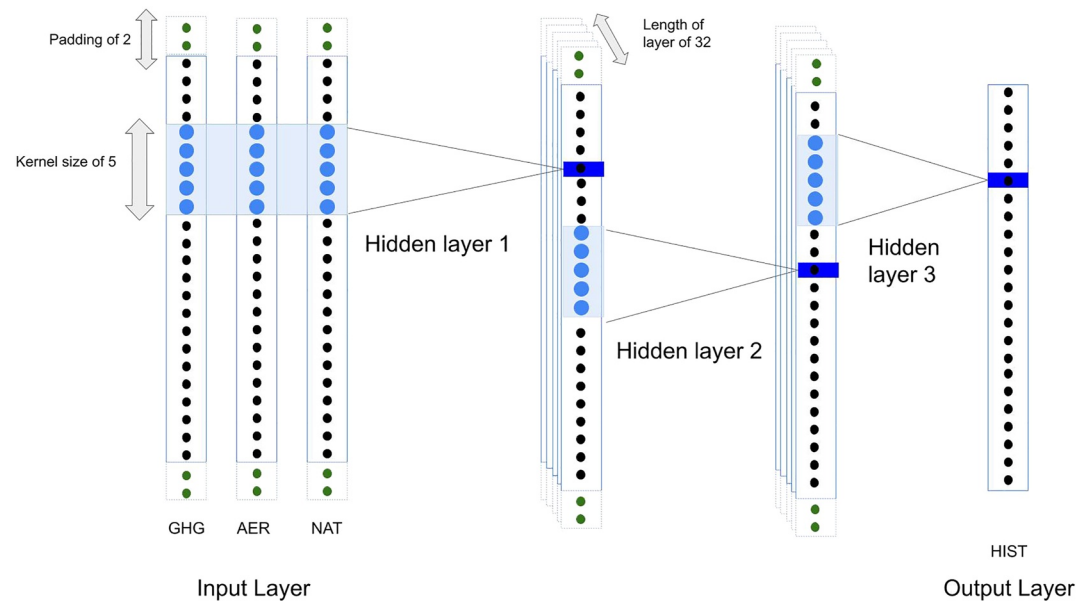


Figure 3. Diagram of the convolutional neural network used. Each white-filled blue rectangles represents a time series of 115 years. The input layer is shown on the left, the hidden layers in the middle and the output layer on the right. Light blue-filled rectangles represent the kernels of the different hidden layers. Dark blue-filled rectangles represent the output of the kernel. Zero-padding is shown in green with dotted lines.

Three steps are applied. First, a climate model is randomly selected. Second, we randomly select one 4-tuple from the chosen climate model. Then the CNN is trained using the three GSAT time series dedicated to (GHG, AER, NAT) as input and the GSAT dedicated to HIST as the target. We iterate this process by repeating it $5 \cdot 10^6$ times. A lower number of iterations was found to degrade the backward optimization results (not shown), but the results are otherwise similar when increasing the number of iterations.

A neural network uses hyperparameters which are the variables that determine the network structure and those which determine how the network is trained. The hyper-parameters are chosen using a cross-validation, as detailed in Text S2 and Figure S1 in Supporting Information S1. The chosen architecture has three convolutional hidden layers, a kernel size of 5 for all layers and 32 filters for each layer.

3.2. Performance of the CNN

Before presenting the neural network dedicated to the attribution method in the next section, we investigate the performance of the CNN in estimating the total effect of forcing from the effect of each forcing separately. First, we train the CNN using the data from all models and estimate the mean training root mean square error (RMSE) made in predicting the data for each model separately. Second, we successively train the CNN leaving out the data from one model and estimate the mean cross-validation RMSE in predicting the left-out model data. Because internal variability is included in the training data, we expect the RMSE to exceed the internal variability in all climate models. The training RMSE is within 0.10°C and 0.25°C for the different climate models. Indeed, the models with large training RMSE (Figure 4 blue bars) corresponds to those simulating a large internal variability, as estimated by the standard deviation of the GSAT of the PI simulation (Table 1), where the forced signal is absent.

The CNN also should produce an estimated GSAT similar to the mean output from the training data, which is expected to be similar to the MMM from

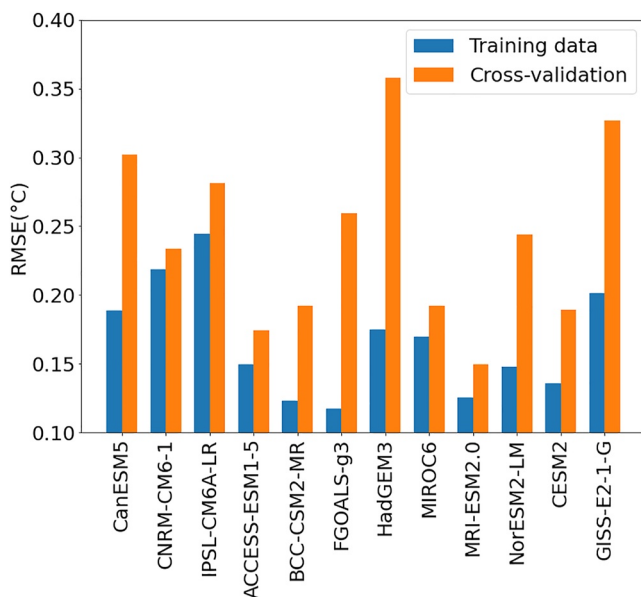


Figure 4. Root mean square error between the convolutional neural network output and the global mean surface air temperature of HIST, in $^\circ\text{C}$, when using (blue bar) the training data and (cross validation, orange bar) when using the data of a model left out in the training.

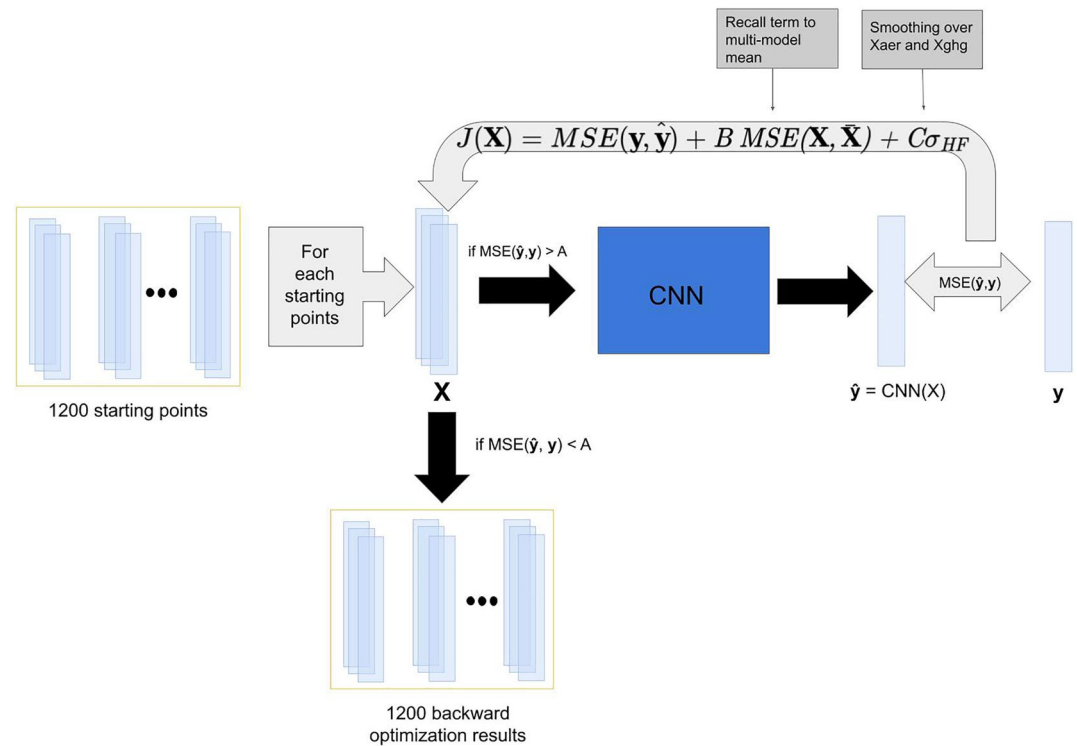


Figure 5. Schematic of the backward optimization attribution process with one entry denoted \mathbf{y} at the right. The 1,200 backward optimization results are at the bottom. The learned convolutional neural network (CNN) is in the middle in dark blue. $J(\mathbf{X})$, the cost function of the backward optimization is on the top. \mathbf{X} denote the optimized input and $\hat{\mathbf{y}}$ denotes its image by the CNN. The 1,200 starting points are on the left.

HIST. The training RMSE may also reflect a forced signal in the HIST simulations distinct from the other models. The amplitude of the RMSE increases to 0.15°C – 0.35°C when using cross-validation. This suggests that the CNN does not overfit. HadGEM3 and, to a lesser extent, FGOALS-g3 and GISS-E2-1-G, show differences much larger than the training RMSE when the data from these models is used for the validation. This might reflect important singularities for these three models, which is probably linked to their singular response to forcings. This might be linked to the equilibrium climate sensitivity, which quantifies the ability of a model to warm up when greenhouse gases increase. It depends on the feedbacks acting in the climate system, and remains poorly constrained by observations (Sherwood et al., 2020). GISS-E2-1-G simulates one of the lowest equilibrium climate sensitivity, while HadGEM2 has one of the highest sensitivities. In addition, FGOALS-g3 simulates almost no response to anthropogenic aerosols (see Figure 1).

3.2.1. Backward Optimization

In this section, we describe how we use the CNN to perform climate change attribution. The backward optimization is a method that infers the most likely input of the CNN from a given output. To attribute climate change from the CNN, we calculate such input, which provides the GSAT attributed to the three forcings from the total GSAT anomaly observed or simulated. This is a neural network interpretation method (Gagne et al., 2019; McGovern et al., 2019; Toms et al., 2020) also known as variational inversion when applied to a geophysical model (Brajjard et al., 2012). A scheme of the procedure is given in Figure 5. This optimal input is determined by minimizing a dedicated cost function and using the backpropagation. The cost function, called J , is:

$$J(\mathbf{X}) = \text{MSE}(\mathbf{y}, \text{CNN}(\mathbf{X})) + B \text{MSE}(\mathbf{X}, \bar{\mathbf{X}}) + C \sigma_{HF} \quad (2)$$

where $\mathbf{X} = (x_{GHG}, x_{AER}, x_{NAT})$ is the optimal input to be determined, that is, a triple of 115-year time series corresponding to the GSAT induced by greenhouse gases, anthropogenic aerosols and natural forcing. $\bar{\mathbf{X}}$ is the three time series obtained with the MMM of the simulations GHG, AER and NAT (see Figure 1, lower panel). MSE denotes the mean squared error. \mathbf{y} is the desired output of the neural network. σ_{HF} is the sum of the time standard

deviation of the high-pass filtered time series obtained from x_{GHG} and x_{AER} using a Lanczos high-pass filter with a window of size 21, a cutoff period of 10 years. B and C are two adjustable real parameters.

The first term on the right-hand side of Equation 2 measures the mean square error between the desired output and the CNN output. The second term, also known as a background term, is applied so that the results are similar to a first guess, taken from the MMM to avoid absurd and nonphysical solutions. Although this term is not standard for the backward optimization of a neural network, it is, however, used for the variational inversion procedure used in data assimilation (Brajaud et al., 2012; Fablet et al., 2021). The last term is used to build smooth GSAT time series for the forcings associated with greenhouse gases and anthropogenic aerosols. Again, this term is not used for the natural forcings, so that the effects from volcanic aerosols remain unsmoothed, with cooling peaks lasting 2–5 years, as expected.

When estimating the optimal input, the initial input is iteratively updated using a back-propagation to minimize $J(\mathbf{X})$ until it is smaller than a fixed value, called A . To reduce the computational cost, the minimization process is stopped after 500 iterations if $J(\mathbf{X})$ does not converge. The backward optimization of a neural network has multiple solutions and the method is sensitive to the initial value used for \mathbf{X} . Therefore, for each of the 12 climate models, we randomly select with repetition 100 (10 during the perfect model approach) triples of the GSAT time series among the members of GHG, AER, and NAT as first guess for the initial states. These initial states are chosen as they represent physically coherent inputs. This provides 1,200 initial physically coherent values for \mathbf{X} which sample the internal climate variability and the spread among the different models. This generates 1,200 backward optimizations. This estimation is empirical and does not account for the internal variability of the target of the backward optimization. For each year, the 90% confidence intervals of the optimal input are then estimated using ± 1.64 standard deviations among the backward optimization results, assuming a Gaussian distribution.

The choice of A (iteration stop threshold), B (background term) and C (smoothing term) was fixed empirically as the other hyperparameters of the neural network. We found that these parameters do not significantly modify the results of the backward optimization (see Text S3, Tables S1 and S2 in Supporting Information S1). We select $A = 0.05$, $B = 0.01$ and $C = 0.1$.

3.3. Regularized Optimal Fingerprints

We evaluate the performance of the neural network based method for detection and attribution by comparing its results to those obtained with the ROF (Ribes et al., 2013). This last method is widely used and has already been applied to the air surface temperature using CMIP6 data by Gillett et al. (2021).

The ROF method is based on a multivariate linear regression and on the assumption that the observed change can be obtained with the sum of the forced anomalies for each forcing (the so-called fingerprints) plus internal variability.

The observed GSAT denoted \mathbf{y} , is given by:

$$\mathbf{y} = \beta \mathbf{X} + \epsilon \quad (3)$$

with $\beta = (\beta_{GHG}, \beta_{AER}, \beta_{NAT})$ the scaling factors and $\mathbf{X} = (X_{GHG}, X_{AER}, X_{NAT})$ the effects of all the forcings on the GSAT. ϵ represents the effect of internal variability, assumed to be a Gaussian white noise.

We use greenhouse gases, anthropogenic aerosols and natural forcings as three individual forcings and neglect the other forcings. \mathbf{X} is estimated in this case by using the MMM of GHG, AER and NAT simulations.

To perform such a regression, a common method is to reduce the dimension of data using the leading empirical orthogonal functions calculated in PI. This reduces the number of spatial dimensions and allows an accurate estimation of the internal variability covariance matrix. But such a method involves an arbitrary choice of the number of EOFs used to truncate the data. The ROF method (Ribes et al., 2013) avoids this arbitrary choice using a regularized estimation of the covariance matrix to estimate the scaling factors.

The response of climate to the i th forcing is detected if β_i is significantly different from zero. If the confidence interval of β_i includes one, this shows consistency between observations and simulated climate model responses. We use the total least square (TLS) method (Allen & Stott, 2003) to perform the regression and estimate the scaling factors, which accounts for the residual internal variability in the MMM. The internal variability is

assumed to be the same in GHG, AER, and NAT members, which prevents the use of different smoothing to the GSAT simulated in GHG and AER, as done for the backward optimization, or in NAT. As the internal variability is largely reduced by the ensemble averaging in the MMM, we estimate the attributable warming in GSAT by $\beta_i X_i$ for the i th forcing. This should lead to an attributable warming similar to $\beta_i \hat{X}_i$ using the estimated X_i by the TLS instead of X_i . Estimates of attributable warming in GSAT for each year can then be obtained by $\sum \beta_i X_i$. Following Gillett et al. (2021), the internal variability is sampled by concatenating all available simulations after subtraction of the mean of the corresponding model ensemble. To account for the subtraction of the ensemble mean, we multiply for each model, the anomalies by $\sqrt{\frac{n}{n-1}}$, where n is the ensemble size. For each simulation, the equivalent size corresponding to the MMM is estimated using:

$$N = \frac{M^2}{\sum_{i=1}^M \frac{1}{n_i}} \quad (4)$$

with M the number of different climate models used (in our Case 12) and n_i the number of members available for the i th climate model. To estimate the uncertainty in the GSAT effect attributable to the i th forcing, it is necessary to take into account the uncertainty of β_i and the internal variability contained in X_i . For each year and forcing, the uncertainty in the attributable GSAT is calculated using 1,000 random draws assuming a Gaussian distribution for both β_i and X_i . The mean and standard errors of β_i are estimated as in Allen and Stott (2003). The mean and standard deviation of X_i are estimated from the size N of the MMM and the standard deviation of the GSAT obtained from the PI runs. We first calculate the standard deviation for each model (as given in Table 1), average the values obtained across models, and then divide by the square root of N . This procedure is valid under the conditions that the uncertainties of β_i and X_i are Gaussian, uncorrelated and small compared to their respective means. The latter hypothesis is not verified for GSAT anomalies close to zero for X_i , such as those obtained in the first decades of our time series (see Figure 1), or for the GSAT of NAT. Thus the uncertainties for the attributable GSAT are to be taken with caution.

4. Attribution Performances

4.1. Performance on Synthetic Data

To investigate the performance of the backward optimization and ROF in the case of non-additive data, we applied the two attributions methods to the synthetic data presented in Section 2.4. Figures 6a and 6c shows the time series of the estimated effect of the three synthetic forcings and f_1 , f_2 , and f_3 the ground truth time series. We use the 100 surrogate time series generated for each forcings and their response denoted F_1 , F_2 , F_3 , and R , instead of the simulated GSAT from GHG, AER, NAT, and HIST, respectively. The 10 R time series remaining are used as pseudo-observation, noted y previously. For the backward optimization, the estimated forced effect f_1 (Figure 6c, red lines) show some variability but is centered around the true f_1 (purple line). For ROF (Figure 6a, red lines), the estimated f_1 are systematically larger than the true f_1 at the end of the time series. Similarly, f_2 (Figure 6c, blue and dark blue lines) is well estimated by the backward optimization, while ROF (Figure 6a) produces an estimated f_2 with an important variability and an overestimation in most of the cases. The f_3 forcing is well estimated by both methods, but with more variability for backward optimization.

The RMSE between the effect estimated by the different attribution methods and the ground truth are shown in 6bd in the form of boxplot. For ROF the mean RMSE value is 0.14°C for f_1 , 0.12°C for f_2 and 0.01°C for f_3 . These values are for backward optimization of 0.05°C for f_1 , 0.04°C for f_2 and 0.04°C for f_3 . Backward optimization therefore provides errors smaller than ROF in case of the non-additive forcing generated, while the use of ROF lead to important errors.

4.2. Evaluation of the Performances in Attributing Climate Changes: Perfect Model Approach

To evaluate the performance of the backward optimization and ROF we use a perfect model approach that relies on climate model data only. This approach consists of using the data from all but one of the climate models to perform our two attribution methods. In the case of backward optimization, this implies that we do not use the data from a climate model during the CNN training phase, in the starting points, or in the MMM calculation. For ROF, the data of a model are not used to construct the climate noise estimate or included in the MMM. We use a HIST member of the test data set (see Section 3.3) from each climate model as the target for the attribution

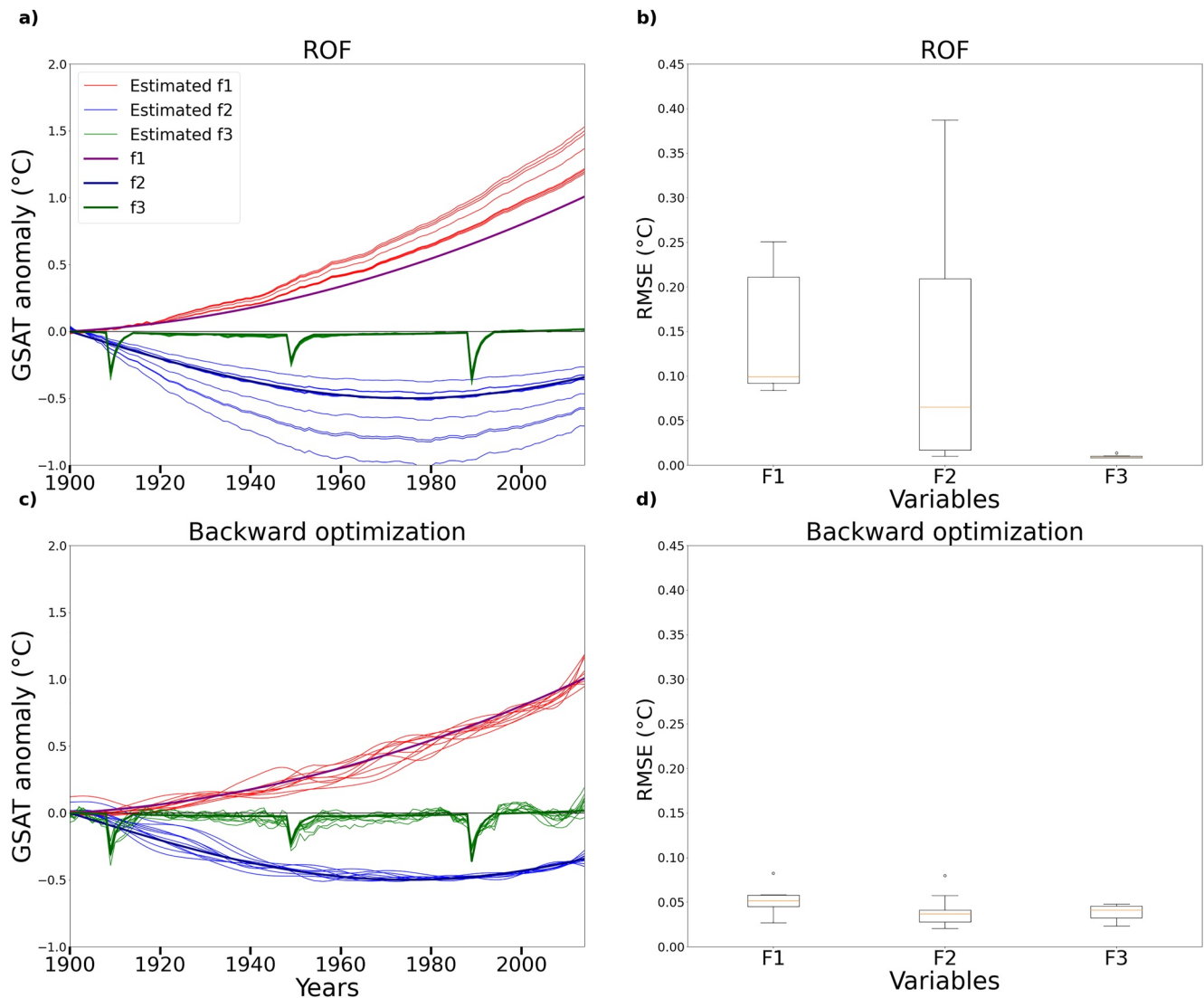


Figure 6. Estimated f_1 , f_2 , and f_3 given by (a) regularized optimal fingerprinting (ROF) and (c) backward optimization. The original f_1 , f_2 , and f_3 ground truth lines are shown in bold. Histograms show the distribution of the root mean square error of the results of (b) ROF and (d) backward optimization compared to the ground truth.

methods. The attributable anomalies associated with each forcing are then compared with the ensemble mean of the GHG, AER and NAT simulations of the removed climate model, even if it includes some residual internal variability, especially when the number of members is small. We use the paradigm that “climate models are statistically indistinguishable from the truth” (Hargreaves, 2010; Ribes et al., 2017; van Oldenborgh et al., 2013), where the difference between observations and models is assumed to be distributed as the difference between any pairs of climate models. We therefore assess the capability of the attribution methods when using observations by investigating only climate models. This approach is called a perfect model approach by analogy with the methods developed for seasonal (Doblas-Reyes et al., 2013) or decadal (Hawkins et al., 2011) climate forecast.

Figure 7 illustrates the attributable anomalies calculated from an HIST member for each climate model. The ensemble means of GHG, AER, and NAT simulations for that climate model are shown for comparison. The differences between the attributable anomalies and the ensemble means of GHG, AER and NAT are also quantified in Figures 8a, 8b, 8c, and 8d with the RMSE and the time mean difference between the two time series. Lastly, the widths of the 90% confidence intervals in 2000–2014 are compared in Figures 8e and 8f.

The two methodologies show a monotonic warming induced by the greenhouse gases that intensified in the 1970s for all climate models. The cooling effect of anthropogenic aerosols is also consistent for both methods, with an intensified cooling in the 1970s, also known as global dimming (Wild, 2009), followed by a stabilization in

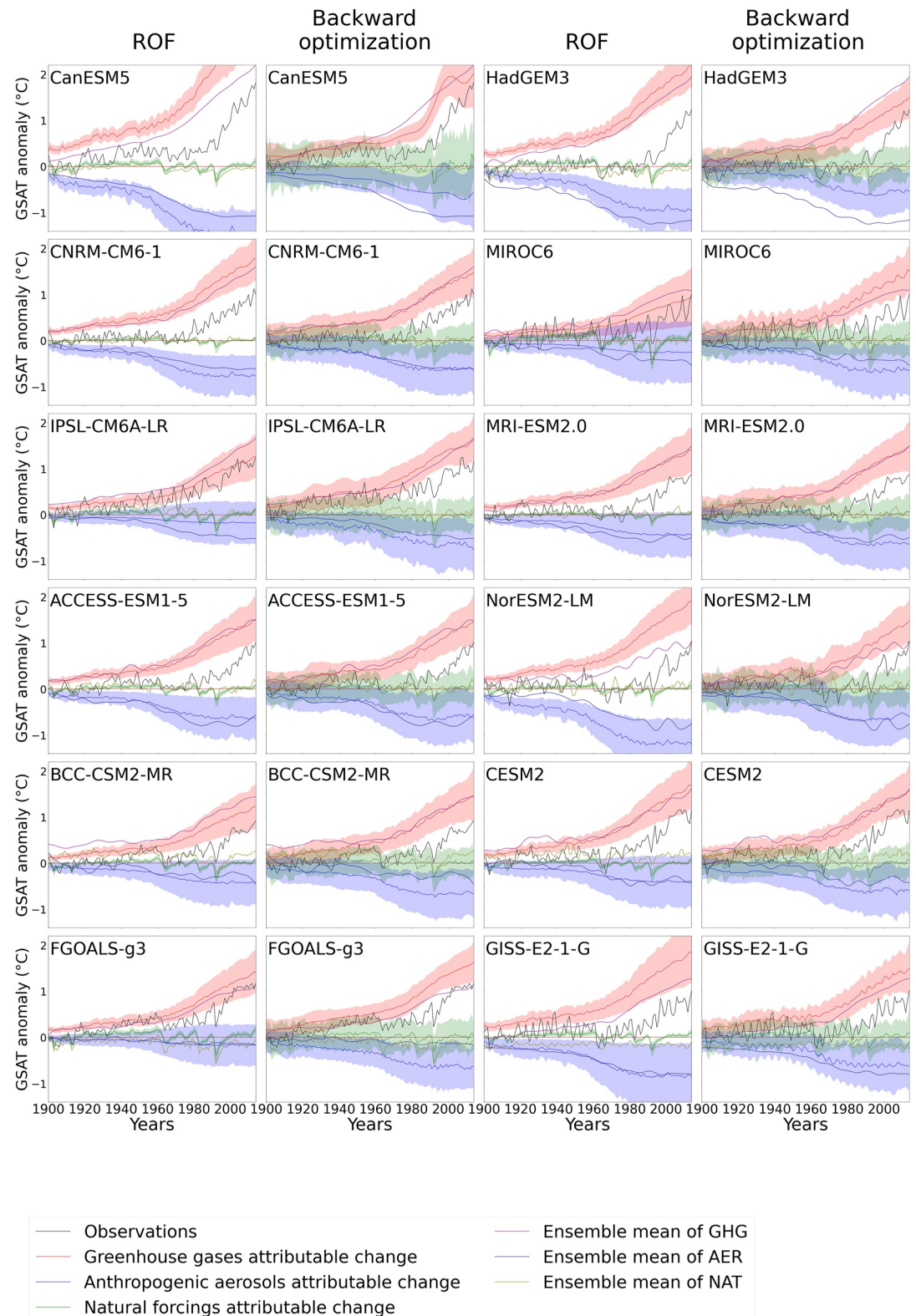


Figure 7. Attributable global mean surface air temperature (GSAT) in $^{\circ}\text{C}$ calculated for regularized optimal fingerprinting (ROF) and backward optimization from (black line) a HIST member. The GSAT is decomposed into the attributable changes due to (red line) greenhouse gases; (blue line) anthropogenic aerosols and (green line) natural forcings. For comparison, the ensemble mean of (purple line) GHG, (dark blue line) AER and (beige line) NAT is indicated. Color shades show the 90% confidence intervals of the attributed GSAT.

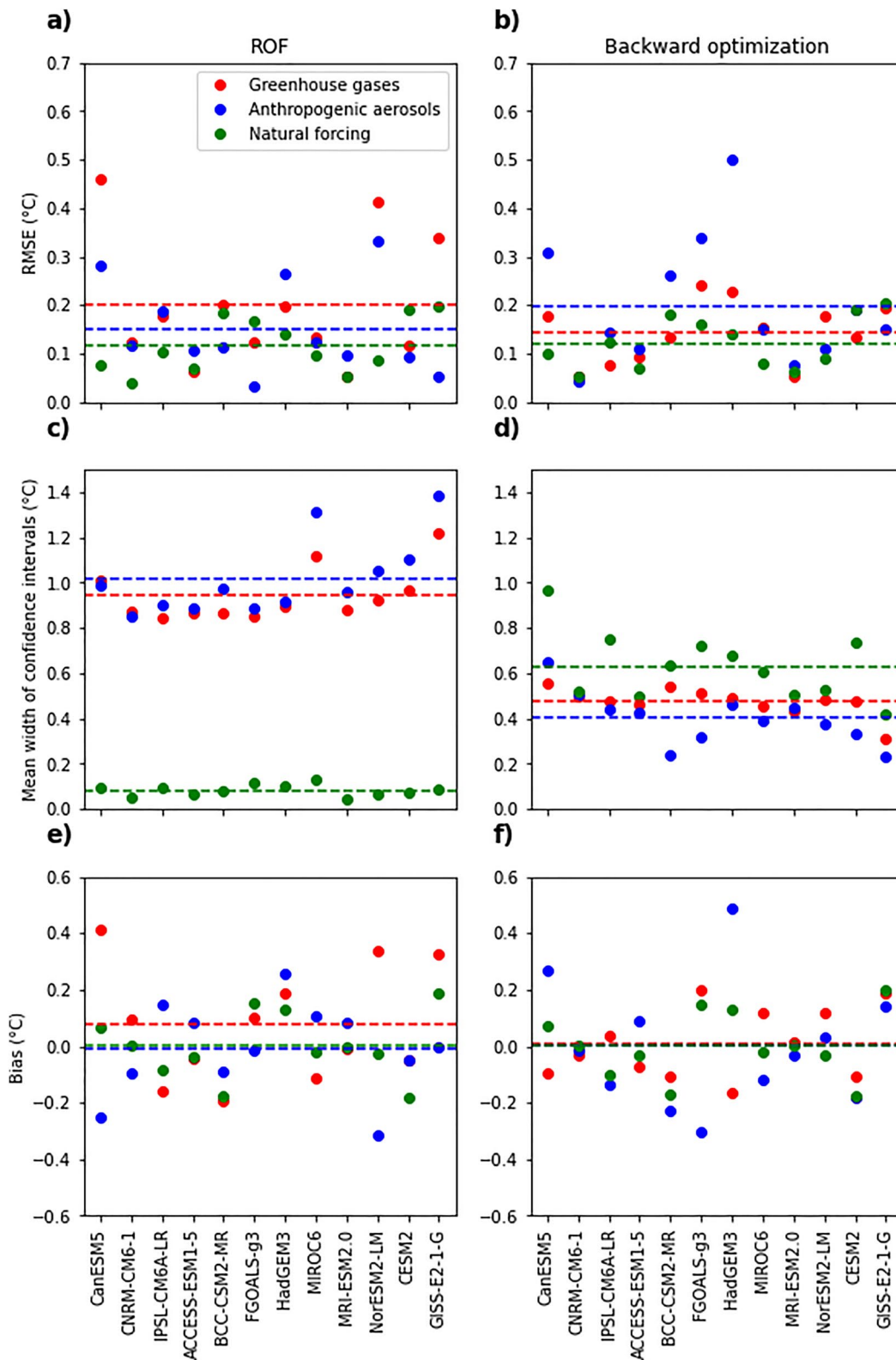


Figure 8.

the 2000s. Lastly, the changes attributable to natural forcings are small in both methods, except for the cooling following the major volcanic eruptions.

For the backward optimization, the RMSE is 0.14°C, 0.20°C, and 0.12°C when averaged across the 12 models for the effects of greenhouse gases, anthropogenic aerosols and natural forcing, respectively (see dashed line in Figure 8a). ROF provides an average RMSE of 0.20°C, 0.15°C, and 0.12°C for these forcings (dashed lines in Figure 8b), so the errors are similar in both methods. Moreover, ROF shows an average positive bias of 0.09°C for greenhouse gases. All other biases for ROF and for backward optimization are almost zero. ROF, therefore seems to overestimate the effect of greenhouse gases which is not the case of the backward optimization.

However, RMSE and biases are affected by the residual internal variability included in ensemble means especially when only a few members are available. The RMSE and biases are therefore weak indicators for models with few members. The width of the confidence intervals for greenhouse gases and anthropogenic aerosols obtained with the backward optimization are smaller than those obtained with ROF from the 1970s, while they are larger from 1900 to 1940. Although the uncertainty provided by the confidence intervals of ROF was verified using a perfect model approach in Gillett et al. (2021), some authors suggested that ROF underestimates such uncertainty because of insufficient consideration in the internal variability (DelSole et al., 2019; C. Li et al., 2021). This suggests that the confidence intervals given by the backward optimization are also underestimated, and that further improvements would be needed to evaluate them in more details.

The width of the confidence intervals for the effect of natural forcing (Figures 8c and 8d, green points) is in ROF much lower than this obtained with the backward optimization. This might be explained by the calculation of the confidence intervals of ROF which is not adapted to small anomalies (see Section 3.3) as obtained for natural forcings. Moreover, we evaluate the uncertainty for the backward optimization by sampling both the inter-model and internal variability contained in the starting points, so that the confidence intervals are rather homogeneous in time and for the three forcings. We suggest that both estimations need to be refined using larger ensembles of simulations. This would allow a more systematic assessment of the uncertainties using the perfect model approach.

Figures 7 and 8 also show that the cooling from anthropogenic aerosols is overestimated in FGOALS-g3 in backward optimization results compared to the ensemble mean, and underestimated in CanESM5 and HadGEM3. It is likely that effect of external forcings in these three models is very different from the other models. For instance, FGOALS-g3 simulates a negligible effect for the aerosols in AER (see Figure 1). CanESM5 and HadGEM3 simulate a warming induced by greenhouse gases (see GHG simulation) larger than the other models, probably associated with the important equilibrium climate sensitivity of these models. The backward optimization fails to reproduce these singular behaviors, being mostly governed by the multi-model consensus. The CNN-based method, that is, the backward optimization, shows results less variable between models than ROF. The backward optimization attributable changes are more consistent with the multi-model consensus, which is hardly affected by removing the data from one climate model. In contrast, in ROF the MMM time series is rescaled with the scaling factors (see Section 3.3). This leads to important errors when the data used as pseudo-observation is taken from a model with a large sensitivity (see for instance CanESM5).

Figure 7 is only based on the use of a single historical simulation for each model. Therefore, we also investigate if the attributable changes are affected by a modification of the historical member. The attributable GSAT is estimated with the two methods from the 10 HIST IPSL-CM6-LR member from the test data (see Section 3.1.1). The RMSEs, the biases and the width of confidence intervals are obtained with respect to the ensemble mean of the single-forcing simulations of the IPSL-CM6-LR model (Figure S2 in Supporting Information S1). Backward optimization presents much less variable results between members than ROF in terms of RMSE or bias, except for natural forcing. The amplitude of the confidence intervals is slightly increases for the backward optimization compared to ROF. It results that backward optimization is less affected by internal variability than ROF.

Figure 8. Performances of attribution methods using a perfect model approach. (a) Root mean square error when using regularized optimal fingerprinting (ROF) for the attributable global mean surface air temperature (GSAT) anomaly of (red) greenhouse gases, (blue) anthropogenic aerosols, (green) natural forcing. (b) Same as (a) for the backward optimization. (c) Width of the 90% percent confidence intervals in 2000–2014 when using ROF. (d) Same as (c) but for backward optimization (e) Time mean difference between the estimated and ensemble mean GSAT attributable to the forcings when using ROF. (f) Same as (e) for backward optimization. Dashed lines shows average values across the 12 climate models.

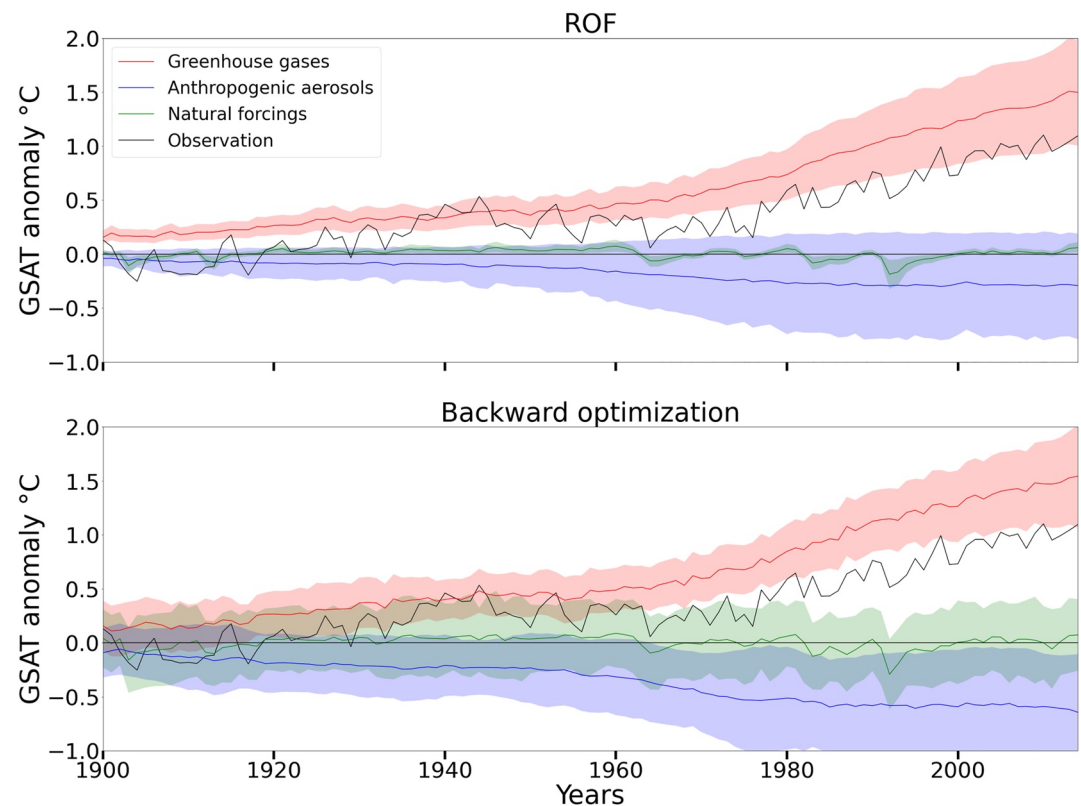


Figure 9. (Top) Attributable global mean surface air temperature (GSAT) anomaly, in $^{\circ}\text{C}$, as given by regularized optimal fingerprinting (ROF) for the effect of the (red) greenhouse gases, (green) natural forcings and (blue) anthropogenic aerosols. The black line shows the observed GSAT. The color shade shows the 90% confidence interval. (Bottom): Same as top, but for backward optimization.

4.3. Attribution of the Observed GSAT

After studying the performance of ROF and backward optimization for synthetic data and in a perfect model approach, we apply both methods to the observed GSAT anomalies.

The attributable GSAT changes are similar for ROF and backward optimization (see Figure 9). For example, in 2000–2014, ROF provides a GSAT attributable to greenhouse gases of 1.28°C (90% confidence interval of $[0.85^{\circ}\text{C}, 1.71^{\circ}\text{C}]$), while it is -0.33°C ($[-0.80^{\circ}\text{C}, 0.12^{\circ}\text{C}]$) for anthropogenic aerosols and 0.01°C ($[0.0^{\circ}\text{C}, 0.02^{\circ}\text{C}]$) for natural forcing. In comparison, backward optimization finds attributable changes of 1.42°C ($[1.03^{\circ}\text{C}, 1.80^{\circ}\text{C}]$), -0.61°C ($[-1.16^{\circ}\text{C}, -0.06^{\circ}\text{C}]$) and 0.02°C ($[-0.33^{\circ}\text{C}, 0.38^{\circ}\text{C}]$), respectively, for these three forcings. Nevertheless, backward optimization provides more noisy time series and more cooling during volcanic eruptions. The similarity of the results between ROF and backward optimization suggests that the GSAT changes are largely additive as found in Marvel et al. (2015) or Shiogama et al. (2013).

The attributable changes of the GSAT given by ROF are much comparable to that of Gillett et al. (2021) who studied the effect of other forcings (land use and ozone) together with the greenhouse gases. Their results for the 2010–2019 decade provide a 5%–95% range of the attributable warming of $[1.2^{\circ}\text{C}, 1.9^{\circ}\text{C}]$ for greenhouse gases and other forcings, $[-0.7^{\circ}\text{C}, -0.1^{\circ}\text{C}]$ for anthropogenic aerosols and $[0.01^{\circ}\text{C}, 0.06^{\circ}\text{C}]$ for natural forcing. We verified that the ROF results shown in Figure 9 remain similar when we take into account other forcings together with the greenhouse gases influence (see Figure S3 in Supporting Information S1).

Backward optimization shows a slightly smaller uncertainty for greenhouse gases and anthropogenic aerosols than ROF toward the end of the time series, but a larger uncertainty range for natural forcings, as found and discussed in Section 4.2. We can note that the reconstruction of the observations by the backward optimization is by construction very close to the observations (see Figure S4 in Supporting Information S1) and captures most of the internal variability contained within the observations.

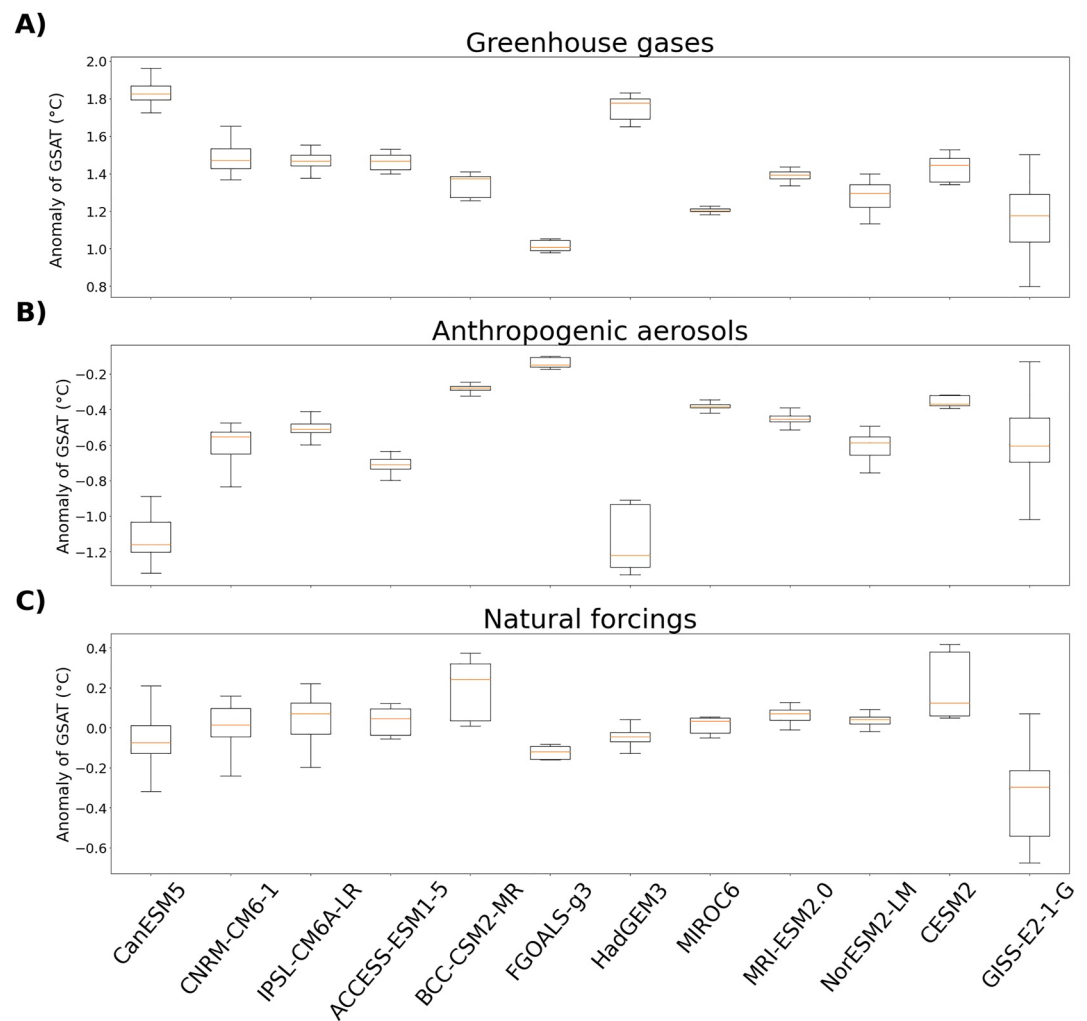


Figure 10. Boxplots of the attributable changes in 2000–2014 when using observation and backward optimization, classified according to the climate model used as initial inputs for (a) the greenhouse gases (b) the anthropogenic aerosols and (c) the natural forcings.

4.4. Focus on the Main Backward Optimization Results

The backward optimization uncertainties are computed sampling various initial inputs. Backward optimization is often used with an all-zeros starting point (Toms et al., 2020) even if McGovern et al. (2019) have optimized the initial inputs by using coherent starting points as done in the present study. Figure 10 shows the boxplots of the attributable changes in 2000–2014 when using the observations and backward optimization, as previously discussed in Section 4.3, classified according to the climate models used for the initial input.

The attributable changes produced by the backward optimization are influenced by the climate model used to generate the initial input. For example, CanESM5 simulates large warming in response to greenhouse gases (see Figure 1), probably linked to its large equilibrium climate sensitivity. When using the outputs of CanESM5 as initial input of the backward optimization, large attributable changes are obtained for both the greenhouse gases and the anthropogenic aerosols. On the other hand, when using an initial input from FGOALS-g3 the changes due to the greenhouse gases and the anthropogenic aerosols are small. For each forcing, we analyze the dispersion of the GSAT anomalies over the years 2000–2014 by estimating the mean GSAT attributable to the use of all starting points for each of the 12 climate models. The variability explained by the model is calculated by is the standard deviation across these 12 attributable GSAT. The residual variability which accounts for the internal variability of the starting points is estimated by the standard deviation of the 1,200 attributable GSAT after subtracting for each time series the average response obtained with their respective climate model. The standard deviation explained

by the model of the starting point is 0.22°C for greenhouse gases, 0.29°C for anthropogenic aerosols and 0.14°C for natural forcings. The residual standard deviation is of 0.06°C for the greenhouse gases, 0.09°C for the anthropogenic aerosols and 0.1°C for the natural forcings. The residual variance therefore is smaller than this associated with the climate model for each forcing, especially for the greenhouse gases. The range of attribution results is about 1°C for all forcings, with some particular models providing attributable anomalies at the head or the tail of the inter-models distribution when used as starting point. Removing or modifying these outliers to improve the backward optimization results have been considered. However, selecting these initial inputs may imply a selection of climate models which needs to be associated with a careful investigation of the physical mechanisms (Coquard et al., 2004).

5. Discussion and Conclusion

We present a method for detection and attribution of climate data based on a backward optimization of a CNN. We trained the CNN on the simulated GSATs obtained from outputs of 12 CMIP6 climate models. We then performed a backward optimization to estimate the attributable changes. This methodology does not assume that the effects of the external forcings are additive. Such additivity implies that the total changes simulated by the forcings can be obtained by the sum of the changes due to the individual forcings. The additivity assumption is an important limitation when focusing on precipitation (Marvel et al., 2015) or at regional scale (Deng et al., 2020; Pope et al., 2020). We evaluated the effect of internal variability and model dispersion by using different starting points sampling the simulated distributions. We compared the results of the CNN backward optimization with those obtained using the ROF (Allen & Stott, 2003; Ribes et al., 2013). In order to assess the ability of backward optimization to deal with non-additivities in forcing compared to ROF we used synthetic data, which, unlike GSAT, have a strong non-additive behavior. In that case, the backward optimization results are more similar to the true forced effect of the forcings than when using ROF which assumes additivity. To see if this result can be generalized additional investigations need to be conducted using either different synthetic data or real non-additive climate data, as for instance the precipitation field.

We also designed a perfect model approach to evaluate the skill of the two methods. We successively removed the data of each climate model and used an historical member of the removed climate model as pseudo-observation. The attributable changes of each forcing are then compared to their actual effect simulated in the corresponding ensemble mean of single-forcing simulations. Backward optimization is found to provide performances similar to that obtained with ROFs in terms of RMSEs or bias. The confidence intervals of the backward optimization are smaller for greenhouse gases and anthropogenic aerosols in the last years of the studied period and much larger for natural forcings than those obtained by ROF. As the calculation of the uncertainty applied in ROF has been previously shown to be also underestimated (DeSole et al., 2019), this suggests that backward optimization leads to an even larger underestimation. This might be linked to the internal variability of the target time series, which is not accounted for in the neural network-based method. A solution to solve this issue would be to generate surrogate time series for the backward optimization and repeat the backward optimization. Larger ensemble of single forcing simulations, such as those proposed in the Large Ensemble Single Forcing Model Intercomparison Project (D. M. Smith et al., 2022), would also be required to refine of the estimated errors. In addition, the changes attributable to natural forcings in the backward optimization have a larger uncertainty than the one of ROF. This is suggested to be an artifact of the estimated uncertainty used, which may be flawed for small changes. Many aspects of the backward optimization can be improved in future works. The backward optimization process can also be improved by giving weights based on the realistic simulation of the interannual to decadal variability. Indeed, the procedure presented here is designed to produce a close agreement between the reconstructed time series and the observations (or pseudo-observations). As shown in Figure S4 in Supporting Information S1, the reconstructed time series, i.e., the image of the CNN using the backward optimization results, closely follow the observations. The CNN might instead be designed to only reproduce the forced component of the anomalies excluding the internal variability unrelated to climate forcings. A better treatment of the initial state could be also investigated, excluding or penalizing the time series used as initial input when inconsistent with observations. In addition, giving different weights to each climate models according to their performance in reproducing observed features could be considered, such as the observed GSAT evolution in Ribes et al. (2021).

Overall, the attributable changes obtained with the backward optimization are consistent with recent attribution results, as reviewed in Eyring et al. (2020a). This confirms the previous detection and attribution results on the GSAT. This study also shows that neural networks can be used to explore the CMIP databases through the backward optimization presented here. Such a method could be deployed on other physical variables, such as

precipitation. It could also easily be applied to spatial average instead of global mean where the non-additivities could be an obstacle. Lastly, a similar method applied on gridded data could also be considered without major modifications given that CNNs can easily process images.

Data Availability Statement

The CMIP6 data is available through the Earth System Grid Federation (Cinquini et al., 2014) and can be accessed through different international nodes. For example, <https://esgf-node.ipsl.upmc.fr/projects/esgf-ipsl/>. Codes used in this article for the backward optimization and the figures are from Bône (2023) software available freely at <https://doi.org/10.5281/zenodo.7248662>. The ROF results have been obtained using the Eyring et al. (2020b) software (version 2.9.0) that can be freely found at <https://github.com/ESMValGroup/ESMValTool/releases/tag/v2.9.0>.

Acknowledgments

We thank three anonymous reviewers for their careful reading and their insightful comments and suggestions. We acknowledge the support of the SCAI doctoral program managed by the ANR with the reference ANR-20-THIA-0003, the support of the EUR IPSL Climate Graduate School project managed by the ANR under the “Investissements d’avenir” programme with the reference ANR-11-IDEX-0004-17-EURE-0006. This work was performed using HPC resources from GENCI-TGCC A0090107403 and A0110107403, and GENCI-IDRIS AD011013295. Guillaume Gastineau was funded by the JPI climate/JPI ocean ROADMAP project (Grant ANR-19-JPOC-003).

References

- Allen, M. R., & Stott, P. A. (2003). Estimating signal amplitudes in optimal fingerprinting, Part I: Theory. *Climate Dynamics*, 21(5), 477–491. <https://doi.org/10.1007/s00382-003-0313-9>
- Allen, M. R., & Tett, S. F. (1999). Checking for model consistency in optimal fingerprinting. *Climate Dynamics*, 15(6), 419–434. <https://doi.org/10.1007/s003820050291>
- Barnes, E. A., Hurrell, J. W., Ebert-Uphoff, I., Anderson, C., & Anderson, D. (2019). Viewing forced climate patterns through an AI lens. *Geophysical Research Letters*, 46(22), 13389–13398. <https://doi.org/10.1029/2019gl084944>
- Bône, C. (2023). Codes for “Detection and attribution of climate change” [Software]. <https://doi.org/10.5281/zenodo.7248662>
- Boucher, O., Servonnat, J., Albright, A. L., Aumont, O., Balkanski, Y., Bastrikov, V., et al. (2020). Presentation and evaluation of the IPSL-CM6A-LR climate model. *Journal of Advances in Modeling Earth Systems*, 12(7), e2019MS002010. <https://doi.org/10.1029/2019ms002010>
- Brajaard, J., Santer, R., Crépon, M., & Thiria, S. (2012). Atmospheric correction of MERIS data for case-2 waters using a neuro-variational inversion. *Remote Sensing of Environment*, 126, 51–61. <https://doi.org/10.1016/j.rse.2012.07.004>
- Burger, W., & Burge, M. J. (2009). *Principles of digital image processing: Core algorithms*. Springer London.
- Cassou, C., Kushnir, Y., Hawkins, E., Pirani, A., Kucharski, F., Kang, I.-S., & Caltabiano, N. (2018). Decadal climate variability and predictability: Challenges and opportunities. *Bulletin of the American Meteorological Society*, 99(3), 479–490. <https://doi.org/10.1175/bams-d-16-0286.1>
- Choudhary, K., DeCost, B., Chen, C., Jain, A., Tavazza, F., Cohn, R., et al. (2022). Recent advances and applications of deep learning methods in materials science. *npj Computational Materials*, 8(1), 1–26. <https://doi.org/10.1038/s41524-022-00734-6>
- Cinquini, L., Crichton, D., Mattmann, C., Harney, J., Shipman, G., Wang, F., et al. (2014). The Earth System Grid Federation: An open infrastructure for access to distributed geospatial data [Dataset]. *Future Generation Computer Systems*, 36, 400–417. <https://doi.org/10.1016/j.future.2013.07.002>
- Coquard, J., Duffy, P., Taylor, K., & Iorio, J. (2004). Present and future surface climate in the western USA as simulated by 15 global climate models. *Climate Dynamics*, 23(5), 455–472. <https://doi.org/10.1007/s00382-004-0437-6>
- Danabasoglu, G., Lamarque, J.-F., Bacmeister, J., Bailey, D., DuVivier, A., Edwards, J., et al. (2020). The Community Earth System Model version 2 (CESM2). *Journal of Advances in Modeling Earth Systems*, 12(2), e2019MS001916. <https://doi.org/10.1029/2019ms001916>
- DelSole, T., Trenary, L., Yan, X., & Tippet, M. K. (2019). Confidence intervals in optimal fingerprinting. *Climate Dynamics*, 52(7), 4111–4126. <https://doi.org/10.1007/s00382-018-4356-3>
- Deng, J., Dai, A., & Xu, H. (2020). Nonlinear climate responses to increasing CO₂ and anthropogenic aerosols simulated by CESM1. *Journal of Climate*, 33(1), 281–301. <https://doi.org/10.1175/jcli-d-19-0195.1>
- Doblas-Reyes, F. J., García-Serrano, J., Lienert, F., Biescas, A. P., & Rodrigues, L. R. (2013). Seasonal climate predictability and forecasting: Status and prospects. *Wiley Interdisciplinary Reviews: Climate Change*, 4(4), 245–268. <https://doi.org/10.1002/wcc.217>
- Eyring, V., Bock, L., Lauer, A., Righi, M., Schlund, M., Andela, B., et al. (2020a). Earth System Model Evaluation Tool (ESMValTool) v2.0—An extended set of large-scale diagnostics for quasi-operational and comprehensive evaluation of Earth system models in CMIP. *Geoscientific Model Development*, 13(7), 3383–3438. <https://doi.org/10.5194/gmd-13-3383-2020>
- Eyring, V., Bock, L., Lauer, A., Righi, M., Schlund, M., Andela, B., et al. (2020b). Earth System Model Evaluation Tool (ESMValTool) v2.0—An extended set of large-scale diagnostics for quasi-operational and comprehensive evaluation of Earth system models in CMIP [Software]. 13(7). Retrieved from <https://github.com/ESMValGroup/ESMValTool/releases/tag/v2.9.0>
- Eyring, V., Bony, S., Meehl, G. A., Senior, C. A., Stevens, B., Stouffer, R. J., & Taylor, K. E. (2016). Overview of the Coupled Model Inter-comparison Project Phase 6 (CMIP6) experimental design and organization. *Geoscientific Model Development*, 9(5), 1937–1958. <https://doi.org/10.5194/gmd-9-1937-2016>
- Eyring, V., Gillett, N., Rao, K. A., Barimalala, R., Parrillo, M. B., Bellouin, N., et al. (2021). Human influence on the climate system. In *Climate change 2021: The physical science basis. Contribution of Working Group I to the sixth assessment report of the Intergovernmental Panel on Climate Change*. Cambridge University Press.
- Fablet, R., Chapron, B., Drumetz, L., Mémin, E., Pannekoucke, O., & Rousseau, F. (2021). Learning variational data assimilation models and solvers. *Journal of Advances in Modeling Earth Systems*, 13(10), e2021MS002572. <https://doi.org/10.1029/2021ms002572>
- Gagne, D. J., II, Haupt, S. E., Nychka, D. W., & Thompson, G. (2019). Interpretable deep learning for spatial analysis of severe hailstorms. *Monthly Weather Review*, 147(8), 2827–2845. <https://doi.org/10.1175/mwr-d-18-0316.1>
- Gillett, N. P., Kirchmeier-Young, M., Ribes, A., Shiogama, H., Hegerl, G. C., Knutti, R., et al. (2021). Constraining human contributions to observed warming since the pre-industrial period. *Nature Climate Change*, 11(3), 207–212. <https://doi.org/10.1038/s41558-020-00965-9>
- Gillett, N. P., Shiogama, H., Funke, B., Hegerl, G., Knutti, R., Matthes, K., et al. (2016). The Detection and Attribution Model Inter-comparison Project (DAMIP v1.0) contribution to CMIP6. *Geoscientific Model Development*, 9(10), 3685–3697. <https://doi.org/10.5194/gmd-9-3685-2016>
- Good, P., Lowe, J. A., Andrews, T., Wiltshire, A., Chadwick, R., Ridley, J. K., et al. (2015). Nonlinear regional warming with increasing CO₂ concentrations. *Nature Climate Change*, 5(2), 138–142. <https://doi.org/10.1038/nclimate2498>

- Goodfellow, I., Bengio, Y., & Courville, A. (2016). *Deep learning*. MIT press.
- Gulev, S., Thorne, P., Ahn, J., Dentener, F., Domingues, C., Gerland, S., & Vose, R. (2021). Changing state of the climate system. In *Climate change 2021: The physical science basis. Contribution of Working Group I to the sixth assessment report of the Intergovernmental Panel on Climate Change*.
- Gupta, A. S., Jourdain, N. C., Brown, J. N., & Monselesan, D. (2013). Climate drift in the CMIP5 models. *Journal of Climate*, 26(21), 8597–8615. <https://doi.org/10.1175/JCLI-D-12-00521.1>
- Ham, Y.-G., Kim, J.-H., & Luo, J.-J. (2019). Deep learning for multi-year ENSO forecasts. *Nature*, 573(7775), 568–572. <https://doi.org/10.1038/s41586-019-1559-7>
- Hargreaves, J. C. (2010). Skill and uncertainty in climate models. *WIREs Climate Change*, 1(4), 556–564. <https://doi.org/10.1002/wcc.58>
- Hasselmann, K. (1993). Optimal fingerprints for the detection of time-dependent climate change. *Journal of Climate*, 6(10), 1957–1971. [https://doi.org/10.1175/1520-0442\(1993\)006<1957:offido>2.0.co;2](https://doi.org/10.1175/1520-0442(1993)006<1957:offido>2.0.co;2)
- Hawkins, E., Robson, J., Sutton, R., Smith, D., & Keenlyside, N. (2011). Evaluating the potential for statistical decadal predictions of sea surface temperatures with a perfect model approach. *Climate Dynamics*, 37(11), 2495–2509. <https://doi.org/10.1007/s00382-011-1023-3>
- Hobbs, W., Palmer, M. D., & Monselesan, D. (2016). An energy conservation analysis of ocean drift in the CMIP5 global coupled models. *Journal of Climate*, 29(5), 1639–1653. <https://doi.org/10.1175/jcli-d-15-0477.1>
- Irving, D., Hobbs, W., Church, J., & Zika, J. (2021). A mass and energy conservation analysis of drift in the CMIP6 ensemble. *Journal of Climate*, 34(8), 3157–3170. <https://doi.org/10.1175/JCLI-D-20-0281.1>
- Kelley, M., Schmidt, G. A., Nazarenko, L. S., Bauer, S. E., Ruedy, R., Russell, G. L., et al. (2020). GISS-E2. 1: Configurations and climatology. *Journal of Advances in Modeling Earth Systems*, 12(8), e2019MS002025. <https://doi.org/10.1029/2019ms002025>
- Kennedy, J. J., Rayner, N. A., Atkinson, C. P., & Killick, R. E. (2019). An ensemble data set of sea surface temperature change from 1850: The Met Office Hadley Centre HadSST.4.0.0.0 data set. *Journal of Geophysical Research: Atmospheres*, 124(14), 7719–7763. <https://doi.org/10.1029/2018JD029867>
- Kiranyaz, S., Avci, O., Abdeljaber, O., Ince, T., Gabbouj, M., & Inman, D. J. (2021). 1D convolutional neural networks and applications: A survey. *Mechanical Systems and Signal Processing*, 151, 107398. <https://doi.org/10.1016/j.ymssp.2020.107398>
- Labe, Z. M., & Barnes, E. A. (2021). Detecting climate signals using explainable AI with single-forcing large ensembles. *Journal of Advances in Modeling Earth Systems*, 13(6), e2021MS002464. <https://doi.org/10.1029/2021ms002464>
- Lam, R., Sanchez-Gonzalez, A., Willson, M., Wirnsberger, P., Fortunato, M., Pritzel, A., et al. (2022). GraphCast: Learning skillful medium-range global weather forecasting. arXiv preprint arXiv:2212.12794.
- Li, C., Zwiers, F., Zhang, X., Li, G., Sun, Y., & Wehner, M. (2021). Changes in annual extremes of daily temperature and precipitation in CMIP6 models. *Journal of Climate*, 34(9), 3441–3460. <https://doi.org/10.1175/JCLI-D-19-1013.1>
- Li, L., Yu, Y., Tang, Y., Lin, P., Xie, J., Song, M., et al. (2020). The flexible global ocean-atmosphere-land system model grid-point version 3 (FGOALS-g3): Description and evaluation. *Journal of Advances in Modeling Earth Systems*, 12(9), e2019MS002012. <https://doi.org/10.1029/2019ms002012>
- Li, Z., Zhang, W., Jin, F.-F., Stuecker, M. F., Sun, C., Levine, A. F., et al. (2020). A robust relationship between multidecadal global warming rate variations and the Atlantic multidecadal variability. *Climate Dynamics*, 55(7), 1945–1959. <https://doi.org/10.1007/s00382-020-05362-8>
- Marvel, K., Schmidt, G. A., Shindell, D., Bonfils, C., LeGrande, A. N., Nazarenko, L., & Tsigaridis, K. (2015). Do responses to different anthropogenic forcings add linearly in climate models? *Environmental Research Letters*, 10(10), 104010. <https://doi.org/10.1088/1748-9326/10/10/104010>
- Masson-Delmotte, V., Zhai, P., Pirani, A., Connors, S., Péan, C., Berger, S., et al. (2021). 2021: Changing state of the climate system. In *Climate change 2021: The physical science basis. Contribution of Working Group I to the sixth assessment report of the Intergovernmental Panel on Climate Change* (pp. 287–422). Cambridge University Press.
- McGovern, A., Lagerquist, R., Gagne, D. J., Jergensen, G. E., Elmore, K. L., Homeyer, C. R., & Smith, T. (2019). Making the black box more transparent: Understanding the physical implications of machine learning. *Bulletin of the American Meteorological Society*, 100(11), 2175–2199. <https://doi.org/10.1175/bams-d-18-0195.1>
- Meehl, G. A., Hu, A., Santer, B. D., & Xie, S.-P. (2016). Contribution of the Interdecadal Pacific Oscillation to twentieth-century global surface temperature trends. *Nature Climate Change*, 6(11), 1005–1008. <https://doi.org/10.1038/nclimate3107>
- 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(3), e2019JD032361. <https://doi.org/10.1029/2019JD032361>
- Neelin, J. D., Battisti, D. S., Hirst, A. C., Jin, F.-F., Wakata, Y., Yamagata, T., & Zebiak, S. E. (1998). ENSO theory. *Journal of Geophysical Research*, 103(C7), 14261–14290. <https://doi.org/10.1029/97jc03424>
- Osborn, T. J., Jones, P. D., Lister, D. H., Morice, C. P., Simpson, I. R., Winn, J. P., et al. (2021). Land surface air temperature variations across the globe updated to 2019: The CRUTEM5 data set. *Journal of Geophysical Research: Atmospheres*, 126(2), e2019JD032352. <https://doi.org/10.1029/2019JD032352>
- O'Shea, K., & Nash, R. (2015). An introduction to convolutional neural networks. arXiv preprint arXiv:1511.08458.
- Pope, J. O., Orr, A., Marshall, G. J., & Abraham, N. L. (2020). Non-additive response of the high-latitude Southern Hemisphere climate to aerosol forcing in a climate model with interactive chemistry. *Atmospheric Science Letters*, 21(12), e1004. <https://doi.org/10.1002/asl.1004>
- Ribes, A., Planton, S., & Terray, L. (2013). Application of regularised optimal fingerprinting to attribution. Part I: Method, properties and idealised analysis. *Climate Dynamics*, 41(11), 2817–2836. <https://doi.org/10.1007/s00382-013-1735-7>
- Ribes, A., Qasmi, S., & Gillett, N. P. (2021). Making climate projections conditional on historical observations. *Science Advances*, 7(4), eabc0671. <https://doi.org/10.1126/sciadv.abc0671>
- Ribes, A., Zwiers, F. W., Azais, J.-M., & Naveau, P. (2017). A new statistical approach to climate change detection and attribution. *Climate Dynamics*, 48(1), 367–386. <https://doi.org/10.1007/s00382-016-3079-6>
- Richardson, M., Cowtan, K., & Millar, R. J. (2018). Global temperature definition affects achievement of long-term climate goals. *Environmental Research Letters*, 13(5), 054004. <https://doi.org/10.1088/1748-9326/aab305>
- Roberts, M. J., Baker, A., Blockley, E. W., Calvert, D., Coward, A., Hewitt, H. T., et al. (2019). Description of the resolution hierarchy of the global coupled HadGEM3-GC3. 1 Model as used in CMIP6 HighResMIP experiments. *Geoscientific Model Development*, 12(12), 4999–5028. <https://doi.org/10.5194/gmd-12-4999-2019>
- Seland, Ø., Bentsen, M., Olivé, D., Toniazzo, T., Gjermundsen, A., Graff, L. S., et al. (2020). Overview of the Norwegian Earth System Model (NorESM2) and key climate response of CMIP6 DECK, historical, and scenario simulations. *Geoscientific Model Development*, 13(12), 6165–6200. <https://doi.org/10.5194/gmd-13-6165-2020>
- Sherwood, S., 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(4), e2019RG000678. <https://doi.org/10.1029/2019rg000678>

- Shiogama, H., Stone, D. A., Nagashima, T., Nozawa, T., & Emori, S. (2013). On the linear additivity of climate forcing-response relationships at global and continental scales. *International Journal of Climatology*, 33(11), 2542–2550. <https://doi.org/10.1002/joc.3607>
- 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, D. M., Gillett, N. P., Simpson, I. R., Athanasiadis, P. J., Baehr, J., Bethke, I., et al. (2022). Attribution of multi-annual to decadal changes in the climate system: The Large Ensemble Single Forcing Model Intercomparison Project (LESFMIP). *Frontiers in Climate*, 4, 955414. <https://doi.org/10.3389/fclim.2022.955414>
- Swart, N. C., Cole, J. N., Kharin, V. V., Lazare, M., Scinocca, J. F., Gillett, N. P., et al. (2019). The Canadian Earth System Model version 5 (CanESM5.0.3). *Geoscientific Model Development*, 12(11), 4823–4873. <https://doi.org/10.5194/gmd-12-4823-2019>
- Tatebe, H., Ogura, T., Nitta, T., Komuro, Y., Ogochi, K., Takemura, T., et al. (2019). Description and basic evaluation of simulated mean state, internal variability, and climate sensitivity in MIROC6. *Geoscientific Model Development*, 12(7), 2727–2765. <https://doi.org/10.5194/gmd-12-2727-2019>
- Toms, B. A., Barnes, E. A., & Ebert-Uphoff, I. (2020). Physically interpretable neural networks for the geosciences: Applications to earth system variability. *Journal of Advances in Modeling Earth Systems*, 12(9), e2019MS002002. <https://doi.org/10.1029/2019ms002002>
- van Oldenborgh, G. J., Reyes, F. J. D., Drijfhout, S. S., & Hawkins, E. (2013). Reliability of regional climate model trends. *Environmental Research Letters*, 8(1), 014055. <https://doi.org/10.1088/1748-9326/8/1/014055>
- Voldoire, A., Saint-Martin, D., Sénési, S., Decharme, B., Alias, A., Chevallier, M., et al. (2019). Evaluation of CMIP6 deck experiments with CNRM-CM6-1. *Journal of Advances in Modeling Earth Systems*, 11(7), 2177–2213. <https://doi.org/10.1029/2019ms001683>
- Wild, M. (2009). Global dimming and brightening: A review. *Journal of Geophysical Research*, 114(D10), D00D16. <https://doi.org/10.1029/2008JD011470>
- Wu, T., Lu, Y., Fang, Y., Xin, X., Li, L., Li, W., et al. (2019). The Beijing Climate Center Climate System Model (BCC-CSM): The main progress from CMIP5 to CMIP6. *Geoscientific Model Development*, 12(4), 1573–1600. <https://doi.org/10.5194/gmd-12-1573-2019>
- Yamashita, R., Nishio, M., Do, R. K. G., & Togashi, K. (2018). Convolutional neural networks: An overview and application in radiology. *Insights Into Imaging*, 9(4), 611–629. <https://doi.org/10.1007/s13244-018-0639-9>
- Yukimoto, S., Kawai, H., Koshiro, T., Oshima, N., Yoshida, K., Urakawa, S., et al. (2019). The Meteorological Research Institute Earth System Model version 2.0, MRI-ESM2.0: Description and basic evaluation of the physical component. *Journal of the Meteorological Society of Japan. Series II*, 97(5), 931–965. <https://doi.org/10.2151/jmsj.2019-051>
- Ziehn, T., Chamberlain, M. A., Law, R. M., Lenton, A., Bodman, R. W., Dix, M., et al. (2020). The Australian Earth System Model: ACCESS-ESM1.5. *Journal of Southern Hemisphere Earth Systems Science*, 70(1), 193–214. <https://doi.org/10.1071/es19035>