Forced Component Estimation Statistical Method Intercomparison Project (ForceSMIP)

3	Robert C. J. Wills ^a , Clara Deser ^b , Karen A. McKinnon ^c , Adam Phillips ^b , Stephen Po-Chedley ^d ,
4	Sebastian Sippel ^e , Anna L. Merrifield ^a , Constantin Bône ^f , Céline Bonfils ^d , Gustau
5	Camps-Valls ^{g} , Stephen Cropper ^{c} , Charlotte Connolly ^{h} , Shiheng Duan ^{d} , Homer Durand ^{g} ,
6	Alexander Feigin ⁱ , M. A. Fernandez ^h , Guillaume Gastineau ^f , Andrei Gavrilov ^{i,g} , Emily
7	Gordon ^j , Moritz Günther ^k , Maren Höver ^{l,a} , Sergey Kravtsov ^m , Yan-Ning Kuo ⁿ , Justin Lien ^o ,
8	Gavin D. Madakumbura ^c , Nathan Mankovich ^g , Matthew Newman ^p , Jamin Rader ^h , Jia-Rui
9	Shi ^q , Sang-Ik Shin ^{p,r} , Gherardo Varando ^s
10	^a ETH Zurich, Zurich, Switzerland
11	^b National Center for Atmospheric Research, Boulder, Colorado
12	^c University of California Los Angeles, Los Angeles, California
13	^d Lawrence Livermore National Laboratory, Livermore, California
14	^e Leipzig University, Leipzig, Germany
15	^f UMR LOCEAN, IPSL, Sorbonne Université, IRD, CNRS, MNHN, Paris, France
16	^g Image Processing Laboratory, University of Valencia, Valencia, Spain
17	^h Colorado State University, Fort Collins, Colorado
18	ⁱ Gaponov-Grekhov Institute of Applied Physics, Russian Academy of Sciences, Nizhny Novgorod,
19	Russia
20	^j Stanford University, Stanford, California
21	^k Max Planck Institute for Meteorology, Hamburg, Germany
22	¹ Oxford University, Oxford, UK
23	^m University of Wisconsin-Milwaukee, Milwaukee, Wisconsin
24	ⁿ Cornell University, Ithaca, New York
25	^o Tohoku University, Sendai, Japan
26	^p NOAA/Physical Sciences Laboratory, Boulder, Colorado

27	^q New York University, New York City, New York
28	^r CIRES, University of Colorado Boulder, Boulder, Colorado
29	^s Image Processing Laboratory, University of Valencia, Valencia, Spain

30 Corresponding author: Robert Jnglin Wills, r.jnglinwills@usys.ethz.ch

ABSTRACT: Anthropogenic climate change is unfolding rapidly, yet its regional manifestation 31 can be obscured by internal variability. A primary goal of climate science is to identify the 32 externally forced climate response from amongst the noise of internal variability. Separating the 33 forced response from internal variability can be addressed in climate models by using a large 34 ensemble to average over different possible realizations of internal variability. However, with 35 only one realization of the real world, it is a major challenge to isolate the forced response in 36 observations. In the Forced Component Estimation Statistical Method Intercomparison Project 37 (ForceSMIP), contributors used existing and newly developed statistical and machine learning 38 methods to estimate the forced response over 1950-2022 within individual realizations of the 39 climate system. Participants used neural networks, linear inverse models, fingerprinting methods, 40 and low-frequency component analysis, among other approaches. These methods were trained on 41 large ensemble from multiple climate models and then applied to observations. Here we evaluate 42 method performance within large ensembles and investigate the estimates of the forced response in 43 observations. Our results show that many different types of methods are skillful for estimating the 44 forced response, though the relative skill of individual methods varies depending on the variable 45 and evaluation metric. Methods with comparable skill in models can give a wide range of estimates 46 of the forced response in observations, illustrating the epistemic uncertainty in forced response 47 estimates. ForceSMIP gives new insights into the forced response in observations, its uncertainty, 48 and methods for its estimation. 49

SIGNIFICANCE STATEMENT: The ForceSMIP project aims to reduce uncertainty in estimates 50 of the climate response to anthropogenic and other external forcing and to evaluate statistical and 51 machine learning methods designed to estimate the forced response from individual realizations of 52 the climate system. New and existing statistical and machine learning methods are evaluated within 53 climate models, for which the forced response is known. Applying these methods to observations 54 gives an estimate of the real-world forced response. The observational forced response estimate 55 agrees with climate models on the large-scale features but also shows discrepancies that give 56 insights into responses that may not be simulated well by climate models. In some regions with 57 large internal variability, such as the North Atlantic ocean, it remains difficult to determine the 58 relative contributions of anthropogenic forcing and internal variability to historical changes. 59

60 1. Introduction

Climate variability and change is composed of forced and unforced components. The forced 61 component of climate change, or forced response, includes all spatiotemporal changes in climate in 62 response to external forcing. Here we consider the net response to forcing from greenhouse gasses, 63 anthropogenic aerosols, land-use change, stratospheric ozone, and natural forcing (e.g., volcanic 64 sulfur emissions and solar variability). The unforced component is due to internal variability of 65 the climate system, for example due to modes of climate variability such as the El Niño-Southern 66 Oscillation (ENSO), Atlantic multi-decadal variability (AMV), and the North Atlantic Oscillation 67 (NAO). In some regions or variables that are prone to large internal variability, the unforced 68 component can be comparable in magnitude to or larger than the forced component, even in multi-69 decadal trends (Deser et al. 2012, 2014; Lehner et al. 2020). Accurate estimation of the forced and 70 unforced components of regional climate change is critical for the attribution of historical climate 71 changes and the characterization and understanding of climate variability and extremes. 72

In climate models, the forced component can be isolated using large ensembles, where the same climate model is run many times with the same forcing but differences in initial conditions, leading to differences in the phasing of internal variability. For a climate measure of interest, the ensemble mean of a large ensemble gives an estimate of the forced response, with larger ensembles needed for variables with lower signal-to-noise ratio (Milinski et al. 2020). Assuming linear additivity of the forced and unforced components, the difference of an individual realization from the ensemble

mean gives the contribution of internal variability. An example is shown for 1980-2022 SST trends 79 from a single member of the ACCESS-ESM1-5 large ensemble in Fig. 1, where the full trend (Fig. 80 1a) is separated into forced and unforced components (Fig. 1b and c, respectively) based on the 81 ensemble mean. Large ensembles are now a widespread tool used for climate change attribution, 82 climate projections, and studies of climate variability and extremes (Deser et al. 2020). However, 83 there is only a single realization of the actual climate system, and it is therefore substantially 84 harder to separate observed climate change into forced and unforced components, which is critical 85 for evaluating climate models and understanding discrepancies between models and observations 86 (Wills et al. 2022; Blackport and Fyfe 2022; Simpson et al. 2025). 87

Individual studies have used one or more statistical methods to estimate the forced response in 106 observations for various applications. For example, separating the forced and unforced component 107 of AMV and the associated Sahel rainfall variability has received particular attention (Ting et al. 108 2009; Booth et al. 2012; Zhang et al. 2013; Frankcombe et al. 2015; Bellucci et al. 2017; Frankignoul 109 et al. 2017; Wills et al. 2020; Qin et al. 2020; Latif et al. 2022; He et al. 2023). By using different 110 methods to estimate the forced response, these studies have reached widely differing conclusions 111 ranging from the AMV is mostly forced (Booth et al. 2012; Wills et al. 2020; He et al. 2023) 112 to the AMV is mostly internal variability (Zhang et al. 2013; Ting et al. 2009; Qin et al. 2020; 113 Latif et al. 2022), although many of these studies acknowledge the uncertainty in this conclusion. 114 There are also a range of conclusions on the forced and unforced contributions to the multi-decadal 115 modulation of the global warming rate (DelSole et al. 2011; Dai et al. 2015; Stolpe et al. 2017; 116 Kravtsov et al. 2018) and multi-decadal changes in the Pacific SST pattern (Olonscheck et al. 2020; 117 Wills et al. 2022; Seager et al. 2022; Rugenstein et al. 2023) and the Aleutian low (Smith et al. 118 2016; Oudar et al. 2018), among other climate indices. All of these questions would benefit from 119 a systematic comparison of methods for estimating the forced response in observations, and this is 120 what the Forced Component Estimation Statistical Method Intercomparison Project (ForceSMIP) 121 aims to do. 122



Illustration of selected methods and how they are evaluated in ForceSMIP using climate model large Fig. 1. 88 ensembles. ForceSMIP participants generated a forced response estimate for each of 10 unlabeled evaluation 89 members. While the forced response estimate includes spatiotemporal variations across 8 variables over 1950-90 2022, here each panel shows 1980-2022 annual-mean SST trends: (a) A single evaluation member (1B) from 91 a large ensemble, which after the submissions was revealed to be from ACCESS-ESM1-5. (b) The "correct 92 answer" is thus estimated from the 40-member ensemble mean of ACCESS-ESM1-5. (c) The internal variability 93 contribution to the trend in (a) is computed as (a) - (b). (d) The TrainingEM method is rescaled from the ensemble 94 mean of the training models and does not use information from ACCESS-ESM1-5 other than the global-mean 95 surface temperature trend. It is a reference method meant to illustrate the forced response that would be estimated 96 from a multi-model ensemble mean. (e)-(i) Forced response estimates from selected ForceSMIP methods, with 97 names and numbers in the titles corresponding to those in Table 1. (j) Taylor diagram showing root mean square 98 error (RMSE) normalized by the root-mean-square-amplitude of the ensemble mean (colors), the root-mean-99 square-amplitude normalized by the root-mean-square-amplitude of the ensemble mean, i.e., $\sigma_i/\sigma_{\text{REF}}$ (black 100 arcs), and the uncentered pattern correlation r_i (black rays). Each method is shown as a symbol with numbers 101 corresponding to those in Table 1; diamonds show methods that use pattern information from the training models; 102 circles show methods that do not. The raw data (a) is shown as a white star, and the dashed white line shows 103 $\delta RMSE_i/RMSE_{RAW} = \delta r_i / r_{RAW}$. The skill metrics are averaged over the 5 "unseen model" evaluation members 104 as explained in the text. 105

TABLE 1. Statistical and machine learning methods for forced response estimation submitted to ForceSMIP Tier 1. Included is information about the institutions involved in developing the methods, a rough categorization of the method type (NN = neural network), whether the method uses pattern information from the training models, whether the method is applied to multiple field variables at once (e.g., using the SST forced response to inform the precipitation forced response), and the number of tunable parameters in the method (i.e., parameters which can be influenced by the training models; reported by the method contributor). Methods are ordered by the number of tunable parameters, and this numbering is used throughout the text and figures.

#	Name	Institution(s)	Type of Method	Pattern Information	Multi-Field	N Parameters
1	RegGMST	NCAR	Regression	ession No		0
2	4th-Order-Polynomial	N/A	Reference No		No	0
3	10yr-Lowpass	N/A	Reference No		No	0
4	LFCA	ETHZ	LFCA No		No	2
5	LFCA-2	ETHZ	LFCA	No	No	2
6	MF-LFCA	ETHZ	LFCA No		Yes	2
7	MF-LFCA-2	ETHZ	LFCA No		Yes	2
8	LIMnMCA	Cornell, Tohoku	LIM	LIM No		2
9	ICA-lowpass	MPI-M	Other No		No	3
10	LIMopt	ETHZ	LIM No		No	3
11	LIMopt-filter	ETHZ	LIM	No	No	4
12	Colored-LIMnMCA	Cornell, Tohoku	LIM	No	Yes	5
13	DMDc	Valencia	LIM	No	No	75
14	GPCA	Valencia	Causal Inference	No	No	88
15	GPCA-DA	Valencia	Causal Inference	No	Yes	89
16	RegGMST-LENSem	NCAR	Regression	No	Yes	876
17	MLR-Forcing	LLNL	Regression	No	Yes	1.1e4
18	SNMP-OF	ETHZ	Fingerprinting	Yes	Yes	1.0e4
19	AllFinger	LLNL, WHOI, UCLA	Fingerprinting	Yes	No	1.0e4
20	MonthFinger	LLNL, WHO, UCLA	Fingerprinting	Yes	No	1.2e5
21	3DUNet-Fingerprinters	UCLA, LLNL, WHOI	NN	Yes	No	5.4e5
22	EOF-SLR	IAP, Milwaukee	Fingerprinting	Yes	No	O(1e6)
23	LDM-SLR	IAP, Milwaukee	Fingerprinting	Yes	No	O(1e6)
24	Anchor-OPLS	Valencia	Regression	Yes	No	2.1e6
25	UNet3D-LOCEAN	LOCEAN	NN	Yes	Yes	2.7e6
26	TrainingEM	N/A	Reference	Yes	Yes	9.1e6
27	RandomForest	UCLA	Random Forest	Yes	No	1.0e7
28	EncoderDecoder	CSU	NN	Yes	No	2.3e7
29	EnsFMP	ETHZ	Fingerprinting	Yes	Yes	4.5e7
30	ANN-Fingerprinters	LLNL	NN	Yes	No	1e16

Large ensembles provide a perfect-model testbed for methods that estimate the forced response 130 from individual ensemble members, because their ensemble mean gives a good estimate of the 131 true forced response in that model. This has been the approach of several previous studies, which 132 have developed statistical or machine learning (StatML) methods to estimate the forced response 133 in single realizations, evaluated them using large ensembles, and then applied them to observations 134 (Deser et al. 2014; Frankcombe et al. 2015; Frankignoul et al. 2017; Sippel et al. 2019; Wills 135 et al. 2020; Bône et al. 2024; Rader et al. 2025). However, these studies have generally focused on 136 one or two methods compared to some simple reference methods, and there has been no broader 137 systematic intercomparison of methods. Furthermore, these studies have primarily targeted surface 138 temperature and/or precipitation, and it is not clear how well the methods used generalize to other 139 climate variables. ForceSMIP aims to systematically compare various StatML methods for forced 140 response estimation across multiple variables in a common framework. Here we both assess which 141 methods are skillful within the large ensemble testbed and assess the spread of estimated forced 142 responses in observations. 143

The rest of the paper is organized as follows. In Section 2, we present the ForceSMIP framework 144 and the climate model large ensemble and observational datasets used. In Section 3, we describe 145 the 30 StatML methods that have been submitted to ForceSMIP. In Section 4, we evaluate the 146 skill of methods for the spatial patterns of long-term trends across multiple variables, grid-scale 147 spatiotemporal variability, and the temporal evolution of selected climate indices. In Section 5, 148 we show examples of the forced responses in observations based on the most skillful methods. 149 Finally, in Section 6, we draw conclusions and discuss implications, potential applications, and 150 future directions. 151

152 2. ForceSMIP Framework and Data

The overarching idea of ForceSMIP is that community contributors develop and train StatML methods to estimate the forced response from single ensemble members and then apply them to model-based evaluation data and observations. The methods are then evaluated based on their forced response estimates in the model-based evaluation data, where the true forced response is known. Finally, the observational forced response estimates can be compared across methods that have proven skillful in the model testbed.

In order to train their methods, contributors were provided with data from 5 climate model large 159 ensembles (Table 2). The identity of these *training models* was revealed to the participants. Data 160 from all ensemble members of the historical and future scenario simulations was provided for 8 161 climate variables over 1850-2100: sea-surface temperature (SST), 2-meter air temperature (T2m), 162 precipitation (PR), sea-level pressure (SLP), monthly-maximum of daily precipitation (monmaxpr), 163 monthly-maximum of daily-maximum temperature (monmaxtasmax), monthly-minimum of daily-164 minimum temperature (monmintasmin), and zonal-mean atmospheric temperature (zmTa). The 165 first four variables were taken from monthly outputs of tos, tas, pr, and psl, respectively, using the 166 naming conventions of CMIP6 (Eyring et al. 2016). The remaining four variables were processed 167 from daily output of pr, tasmax, and tasmin and monthly output of ta, respectively. All variables 168 were interpolated to a common 2.5° grid following Brunner et al. (2020). Four of the variables were 169 then additionally processed with CDO (Schulzweida 2023) commands to make derived variables: 170 daily pr with monmax to make monmaxpr, daily tasmax with monmax to make monmaxtasmax, 171 daily tasmin with monmin to make monmintasmin, and monthly ta with zonmean to make zmTa, 172 where monmax takes a monthly maximum, monmin takes a monthly minimum, and zonmean takes 173 a zonal mean. After this processing, all variables have two spatial dimensions (lat and pressure for 174 zmTa; lat and lon for all others) and monthly time resolution. 175

After developing and training their methods, the contributors submitted: (1) descriptions and 184 basic information about their methods, (2) their method code, and (3) output from application of 185 their method to estimate the forced response across all variables in 10 evaluation members over the 186 period 1950-2022. For the purposes of ForceSMIP, we use a broad definition of the *forced response*: 187 it includes all spatiotemporal variations in the ensemble mean, thus including climate variations 188 due to natural climate forcings (e.g., volcanic eruptions and solar variability) and anthropogenic 189 influences (e.g., anthropogenic emissions of greenhouse gases and aerosols). Contributors therefore 190 had to submit forced response estimates for all variables at monthly time resolution for all points 191 on the 2.5° analysis grid. Nevertheless, much of the discussion in the hackathon that preceded 192 the method submission focused on 1950-2022 trends or 1980-2022 trends, and many participants 193 focused on skill metrics like the pattern correlation and root-mean-square error (RMSE) in long-194 term linear trends, as shown in Figs. 1 and 2. These figures will be discussed in more detail in 195 Section 4, but the overall idea is that by applying StatML to a single ensemble member (for which 196

TABLE 2. Large ensemble and observational data used in ForceSMIP. The first 5 models are the training 176 models and the next 5 models are "unseen models", which are the source of the evaluation members 1B, 1D, 177 1E, 1G, and 1J used for method evaluation in this paper. Evaluation member 1I is the observational data. "Total 178 Members" indicates the number of members used to compute the ensemble mean, with the number in parenthesis 179 indicating the number of future scenario members if it is different than the number of historical simulation 180 members. CESM2 members are those with smoothed biomass burning (Rodgers et al. 2021). Note that due to 181 data problems for zmTa in some members of EC-Earth3, only 13 (51) of the total ensemble members were used 182 to compute the ensemble mean for this variable. 183

Model	Evaluation Member	Total Members	Future Scenario	Reference
CanESM5	1C (r20i1p2f1)	25	SSP585	Swart et al. (2019)
CESM2	1F (LE 1281.019)	50	SSP370	Rodgers et al. (2021)
MIROC6	1H (r11i1p1f1)	50	SSP585	Tatebe et al. (2019)
MIROC-ES2L	N/A	30	SSP245	Hajima et al. (2020)
MPI-ESM1-2-LR	1A (r23i1p1f1)	30	SSP585	Olonscheck et al. (2023)
ACCESS-ESM1-5	1B (r10i1p1f1)	40	SSP585	Ziehn et al. (2020)
EC-Earth3	1D (r6i1p1f1)	18 (58)	SSP585	Wyser et al. (2021)
GFDL-SPEAR-MED	1E (r3i1p1f1)	30	SSP585	Delworth et al. (2020)
IPSL-CM6A-LR	1G (r3i1p1f1)	33 (11)	SSP245	Boucher et al. (2020)
NorCPM1	1J (r4i1p1f1)	30	SSP245	Bethke et al. (2021)
ERA5/ERSST5	11	1	N/A	Hersbach et al. (2020); Huang et al. (2017)

the trends over 1980-2022 are shown in Figs. 1a and 2b), the forced response estimates submitted by ForceSMIP contributors (Figs. 1d-i and 2d-i) should approximate as closely as possible the ensemble mean of the corresponding large ensemble (Figs. 1b and 2b) by removing the internal variability (Figs. 1c and 2c). The 1980-2022 trends shown here are just one way in which the spatiotemporally resolved forced response estimates are evaluated in Section 4.

The *evaluation members* are individual ensemble members of 9 different climate models (Table 201 2) and observational and reanalysis data that has been processed equivalently to the model data. All evaluation members had the metadata removed so that it was not possible to determine which dataset they came from. Only two of the ForceSMIP organizers (C. Deser and A. Phillips) knew the identity of these evaluation members. Of the 9 model-based evaluation members, 5 were from *unseen models* that were not among the training models. The method evaluation in Section 4 will primarily rely on these 5 unseen-model evaluation members. The forced response estimates for



FIG. 2. Same as Fig. 1, but for precipitation.

the evaluation members will be evaluated against the ensemble means computed over all available ensemble members. Note that for two models (EC-Earth3 and IPSL-CM6A-LR), there are a different number of historical and future scenario members, and in these cases the ensemble mean is computed separately in the historical and scenario simulations and then concatenated. The ensemble mean against which methods are evaluated will still have some internal variability in it, which may lead to a slight underestimation of method skill.

Data from observations and reanalysis are included as one of the evaluation members (11), so 215 that the methods can be evaluated and applied to observations in a single round of forced response 216 submissions. This initial round of "Tier 1" ForceSMIP submissions focuses on 1950-2022, which 217 was chosen based on the availability of reanalysis data over this period. As such, all "observational" 218 data in Tier 1 except SST is actually from ERA5 reanalysis (Hersbach et al. 2020). Daily tasmax, 219 tasmin, and pr were computed from ERA5 hourly 2-meter temperature and rainfall data, and the 220 other variables were computed from monthly ERA5 data. SST is from the NOAA Extended 221 Reconstructed SST version 5 (ERSST5; Huang et al. (2017)). Accordingly, the observational 222 forced response estimates from ForceSMIP Tier 1 will be subject to any biases present in the 223 ERA5 and ERSST5 datasets. Subsequent Tiers of ForceSMIP will focus on different time periods, 224 1900-2023 and 1979-2023, on which different sets of observational data are available. 225

While the training data, evaluation methods, and forced response estimates are all available at monthly temporal resolution, all analysis in this paper focuses on annual or seasonal averages (for SST, T2m, PR, SLP, and zmTA), annual maximums, and annual minimums. The annual maximum of monmaxtasmax is called TXx, the annual maximum of monmaxpr is called Rx1day, and the annual minimum of monmintasmin is called TNn, following standard conventions in the study of extreme events (Zhang et al. 2011).

3. Statistical and Machine Learning Methods for Forced Response Estimation

Thirty StatML methods were submitted to this first tier of the ForceSMIP project. They comprise 233 a diverse mix of approaches including linear regression on global-mean temperature or forcing 234 timeseries, low-frequency component analysis (LFCA), linear dynamical mode methods such as 235 linear inverse models (LIMs), linear fingerprinting methods, and neural networks or other machine 236 learning (ML) methods (Table 1). This includes both established methods (e.g., LFCA; Wills et al. 237 (2020), LIMopt; Frankignoul et al. (2017), and regression on global-mean surface temperature; 238 Ting et al. (2009); Deser and Phillips (2021)) and methods newly created for ForceSMIP. The 239 development of many of these methods began at a hackathon held at NCAR and ETH Zurich in 240 August 2023. These methods are briefly summarized in the following subsections, with key details 241 listed in Table 1. In Table 1 and throughout the text, methods are ordered by their number of 242 tunable parameters, which range from 0 to $O(10^7)$ or higher. More detailed information about all 243

methods can be found in the Supplementary Material, and code for all methods can be found at
https://github.com/ForceSMIP/tier1-methods (a persistent identifier will be issued upon
publication).

a. Linear regression on global-mean temperature or forcing timeseries: RegGMST, RegGMST-

248 LENSem, MLR-Forcing

Many studies of internal variability, including ENSO, AMV, and the Pacific Decadal Oscilla-249 tion (PDO), remove anomalies associated with global-mean sea-surface temperature (GMSST) 250 or global-mean surface temperature (GMST) changes when defining indices of this variability 251 (Trenberth and Shea 2006; Ting et al. 2009; Frankignoul et al. 2017; Deser and Phillips 2021). 252 Underlying these approaches is an implicit estimation of the forced response based on GMSST or 253 GMST, under the assumption that those globally aggregated metrics are good proxies of the forced 254 response. In ForceSMIP, two methods, RegGMST and RegGMST-LENSem, estimate the forced 255 response by regressing each field onto a timeseries of GMST and combining that regression pattern 256 with the GMST timeseries. RegGMST uses regression on GMST from the target evaluation mem-257 ber and RegGMST-LENSem uses regression on the ensemble-mean GMST from the 50 CESM2 258 large ensemble members (Deser and Phillips 2023b). 259

A similar approach is to regress each field onto timeseries representing important external forcings or internal variability. The method MLR-forcing uses a multiple-linear-regression approach to regress each field onto regional aerosol forcing timeseries and timeseries representing the response to various forcings (including greenhouse gasses, volcanic emissions, and solar forcing) and on detrended Niño3.4 indices, estimating the forced response as the components associated with the forcing timeseries.

²⁶⁶ b. Low-frequency component analysis: LFCA, LFCA-2, MF-LFCA, MF-LFCA-2, ICA-lowpass

Low-frequency component analysis (LFCA) is a method to objectively identify the slowest evolving spatial patterns in a dataset, using linear discriminant analysis applied to principal components to find patterns that maximize the ratio of low-frequency to total variance (Schneider and Held 270 2001; Wills et al. 2018, 2020). It has been used both to study decadal climate variability (e.g., Wills et al. 2019) and to separate forced and unforced components of climate change (Wills et al.

2020). Its usage as a method to separate forced and unforced components is based on the under-272 standing that the forced response evolves on a longer timescale than most internal variability, i.e., 273 it is using timescale separation to separate forced and unforced components. The application of 274 LFCA in ForceSMIP follows Wills et al. (2020), using a 10-year lowpass filter and including 1 or 2 275 low-frequency patterns in the forced response estimate (methods LFCA and LFCA-2, respectively). 276 Additionally, the methods MF-LFCA and MF-LFCA-2 apply the same method to two variables at 277 a time by combining each field with SST, or in the case of SST, combining it with T2m, with each 278 field normalized by the trace of its covariance matrix. 279

While not a form of LFCA, the ICA-lowpass method uses independent component analysis (Hyvärinen and Oja 2000), which similarly finds linear combinations of a chosen set of principal components that maximize a variance criterion, in this case the statistical independence of the principal components. ICA-lowpass applies independent component analysis to lowpass filtered data and identifies the forced pattern based on its spatial uniformity, under the assumption that the spatial scales of forced climate change are larger than those of internal variability.

286 c. Linear dynamical mode methods: LIMopt, LIMopt-filter, LIMnMCA, Colored-LIMnMCA,

287 DMDc, GPCA, GPCA-DA

Linear dynamical mode methods aim to describe the spatiotemporal variability in a dataset by 288 a set of linear dynamical equations, which determine the evolution of a field from one timestep 289 to the next. The specific case of the Linear Inverse Model (LIM), where the evolution operator is 290 determined from lagged covariance information, is widely used in climate science (Penland and 291 Sardeshmukh 1995; Alexander et al. 2008). The concept of a least damped mode of a LIM was 292 introduced by Penland and Sardeshmukh (1995) and has been used to separate the ENSO-related 293 or forced variations in a dataset (Compo and Sardeshmukh 2010; Solomon and Newman 2012; 294 Frankignoul et al. 2017; Xu et al. 2022). For ForceSMIP, the LIMopt and LIMopt-filter methods 295 apply the method LIM optimal perturbation pattern and LIM optimal perturbation filter method 296 of Frankignoul et al. (2017) (see also Wills et al. 2020). The LIMnMCA and ColoredLIMnMCA 297 methods combined a similar approach applied to SST with a maximum covariance analysis to find 298 the covariations between SST and the other ForceSMIP variables, an extra step which we will show 299 made it much more successful than other linear dynamical mode methods for non-temperature 300

variables (i.e., PR, SLP, and Rx1day). ColoredLIMnMCA differs from LIMnMCA by the use of a
 LIM for colored Gaussian noise (Lien et al. 2025).

The DMDc is similar in approach to LIMopt, but with a generalization of LIM to include a 303 linear forcing component (Proctor et al. 2016). Similarly, GPCA and GPCA-DA are based on the 304 representation of the data as a combination of an autoregressive process and a forced response, 305 where the forced response is estimated by the "direct Granger effect" of the exogenous forcing 306 signal, and are an extension of the method presented in Varando et al. (2022). Like MLR-Forcing, 307 these methods employ additional forcing timeseries. Compared to GPCA, GPCA-DA additionally 308 uses empirical orthogonal functions (EOFs) of SLP to control against the internal variability they 309 may represent. 310

d. Linear fingerprinting methods: AllFinger, MonthFinger, SNMP-OF, EOF-SLR, LDM-SLR, Anchor-OPLS, EnsFMP

Broadly speaking, linear fingerprinting methods use model-based forced response patterns as 313 an initial guess of the forced response and then estimate the contribution of this pattern to the 314 observations (or an individual ensemble member treated like observations). While traditional 315 uses of fingerprinting for detection and attribution generally aim to find a timeseries indicating 316 the amplitude of the forced response pattern compared to internal variability, the fingerprinting 317 methods in ForceSMIP additionally combine that timeseries with an estimate of the forced pattern. 318 AllFinger and MonthFinger are derived from pattern-based fingerprint analyses (Hasselmann 319 1979; Santer et al. 2023), where the forced pattern fingerprint is obtained by averaging across models 320 and extracting the leading EOF (amplifying the signal and reducing the noise). Observations—or 321 individual model realizations—are projected onto the fingerprint to create a pseudo-PC time 322 series, measuring the similarity between the fingerprint and the target's time-varying patterns. The 323 predicted trend map is reconstructed using the forced pattern fingerprint and the pseudo-PCs. 324

EOF-SLR and LDM-SLR methods first estimate each model's forced response components (timeseries) in a basis of spatial patterns given by either ensemble EOF or linear dynamic mode (LDM) decomposition (Gavrilov et al. 2020, 2024) of multi-model ensemble simulations. Then a set of optimal fingerprinting patterns is trained to deduce the forced response from a single realization in this ensemble. These patterns are constructed to be robust to model uncertainty
 within the training ensemble, and can thus be applied to the unseen data.

Anchor-OPLS is a generalization of the anchor regression framework for fingerprint extraction introduced by Sippel et al. (2021), where forced responses are predicted at every grid point and orthonormalised partial least squares (OPLS) is used instead of ordinary least squares.

SNMP-OF is a combination of signal-to-noise maximizing pattern (SNMP) analysis (Ting et al. 334 2009; Wills et al. 2020) with optimal fingerprinting (Hegerl et al. 1996); it finds SNMPs from the 335 training models and then projects their optimal fingerprint onto observations, finally recomputing 336 a pattern from regression of observations onto the resulting signal-to-noise maximizing timeseries. 337 EnsFMP attempts to combine the two steps into one by applying SNMP analysis to numerous 338 combinations of model ensemble members and observations. Unlike the other fingerprinting 339 methods in ForceSMIP, these two methods recompute a forced response pattern within observations, 340 and they thus stick closer to the raw data. 341

e. Machine learning methods: 3DUNet-Fingerprinters, UNet3D-LOCEAN, RandomForest, EncoderDecoder, ANN-Fingerprinters

ML contributions to ForceSMIP include one based on a recently developed method (UNet3D-344 LOCEAN; Bône et al. 2024), and four methods newly developed for ForceSMIP, including one 345 that has recently been used to attribute the record-high 2023 SST (EncoderDecoder; Rader et al. 346 2025). Architectures used include a type of convolutional neural network called a U-Net (3DUNet-347 Fingerprinters, UNet3D-LOCEAN), Encoder-Decoder neural networks (EncoderDecoder, ANN-348 Fingerprinters), and random forests (RandomForest). Two of the ML methods learn to remove 349 the internal variability (UNet3D-LOCEAN, EncoderDecoder), and the other three learn to esti-350 mate the forced response (3DUNet-Fingerprinters, ANN-Fingerprinters, RandomForest). ANN-351 Fingerprinters additionally uses the year as one of the inputs. The ML methods used in this 352 study vary in complexity (e.g., N Parameters in Table 2) and employ different parameter tuning 353 and training strategies. Interestingly, the U-Nets trained on the internal variability and the forced 354 component exhibit different strengths across variables (Section 4). 355

³⁵⁶ f. Reference methods: 4th-Order-Polynomial, 10yr-Lowpass, TrainingEM

In addition to the methods submitted to ForceSMIP, we compare against 3 reference methods, 357 which involve minimal processing of either the raw data or the training-data ensemble mean. Two 358 of the reference methods are simple methods to remove high-frequency noise in the raw data. 359 4th-Order-Polynomial estimates the forced response as a 4th-order-polynomial fit to timeseries of 360 each variable at each grid point. It has been used to estimate the forced response in a seminal 361 paper by Hawkins and Sutton (2009) and later tested in large ensembles by Lehner et al. (2020). 362 10-yr-Lowpass estimates the forced response as all variability left after application of a 10-yr 363 Lanczos lowpass filter. 364

While the first two reference methods are based entirely on the data within the single realization 365 of interest, the third reference method, TrainingEM, represents an opposite extreme where most 366 information is taken from the training data. TrainingEM simply takes the multi-model ensemble 367 mean of the 5 training models as the forced response estimate and rescales it by a constant so that 368 it has the same GMST trend over 1950-2022 as the single realization of interest. This is similar 369 to the scaling method introduced by Steinman et al. (2015) and evaluated by Frankcombe et al. 370 (2015). TrainingEM thus represents a type of null hypothesis where climate models have a perfect 371 estimate of the forced response, up to a rescaling based on differences in climate sensitivity. 372

4. Method Evaluation

In order to evaluate the skill of the ForceSMIP methods in isolating the forced response in 374 individual realizations of the climate system, we focus on their skill in determining the forced 375 response in the 5 unseen climate models (i.e., those not in the training dataset) from a single 376 member of their large ensembles. However, the results are not systematically different in the 4 377 evaluation members that were part of the training data (Fig. S1). The forced response estimates 378 include monthly values globally for 1950-2022, so there are many metrics on which they could be 379 evaluated. We will focus here on skill in estimating long-term forced trends, the grid-scale temporal 380 evolution of the forced response, and the forced response in an illustrative set of large-scale climate 381 indices. 382

383 a. Long-term trends

Our method for evaluating method skill in isolating the forced component of long-term trends 384 can be visualized in Figs. 1 and 2, showing estimates of forced 1980-2022 annual-mean SST 385 and PR trends from a single evaluation member. The forced trend estimate from each method 386 (panels d-l) is compared against the true forced response, as estimated by the ensemble mean of 387 the corresponding large ensemble (panel b). For comparison, we also show how well the linear 388 trend in the raw data from the evaluation member approximates the true forced response (panel a), 389 which is a reference point we expect methods to improve upon. The difference between the full 390 trend in the raw data and the ensemble-mean forced trend is the contribution of internal variability 391 (panel c), which the methods aim to remove. 392

We quantify the skill of each method's estimate of the forced trend pattern \mathbf{f}_i compared to the true forced trend pattern \mathbf{f}_0 in terms of:

- ³⁹⁵ 1. the uncentered pattern correlation, or cosine similarity, $r_i = \langle \mathbf{f}_i, \mathbf{f}_0 \rangle \|\mathbf{f}_i\|^{-1} \|\mathbf{f}_0\|^{-1}$, where $\langle \cdot, \cdot \rangle$ ³⁹⁶ indicates an area-weighted inner product, $\|\cdot\| = p^{-1}\sqrt{\langle \cdot, \cdot \rangle}$ indicates an area-weighted inner-³⁹⁷ product norm, and *p* is the total number of grid cells,
- 2. RMSE_{*i*} = $\|\mathbf{f}_i \mathbf{f}_0\|$ normalized by the amplitude of the true forced trend pattern $\sigma_0 = \|\mathbf{f}_0\|$, where each method's normalized RMSE is hereafter referred to as nRMSE_{*i*}, and
- 400 3. the amplitude ratio of the predicted and true forced trend patterns (σ_i/σ_0).

The root mean square over the 5 unseen-model evaluation members of each method's nRMSE_{*i*} and forced trend pattern amplitude $\sigma_i = ||\mathbf{f}_i||$ is plotted on a Taylor diagram (Figs. 1j and 2j). The colored shading shows nRMSE_{*i*}, the curved black arcs show contours of the amplitude ratio of the predicted and true forced trend patterns (σ_i/σ_0), and the black rays show contours of the uncentered pattern correlation r_i . Because these three metrics are inter-related, the uncentered pattern correlation r_i shown for each method in the Taylor diagrams is determined from the other two variables by:

$$r_{i} = \frac{\sigma_{i}^{2} + \sigma_{0}^{2} - \text{RMSE}_{i}^{2}}{2\sigma_{i}\sigma_{0}} = \frac{1 + (\sigma_{i}/\sigma_{0})^{2} - \text{nRMSE}_{i}^{2}}{2\sigma_{i}/\sigma_{0}}.$$
 (1)

This equation is exact when applied to a single evaluation member but is approximate when applied to the averages over 5 members in the Taylor diagrams. Note that the Taylor diagrams in this paper do not show the full quadrant, as is traditional (Taylor 2001). Rather, they zoom in on the regions
where the points are. Our variant on the Taylor diagram is partially inspired by the "solar diagram"
of Wadoux et al. (2022), however, in our case the quantitative information remains the same as in
a traditional Taylor diagram.

One noteworthy observation from Figs. 1 and 2 is that methods that do not use pattern information 414 from the training models (methods 1-17; shown with circular symbols in the Taylor diagrams; 415 hereafter simple methods) estimate forced trends that look more like the raw trend from the 416 evaluation member (Fig. 1e-f, cf. Fig. 1a; 2e-f, cf. Fig. 2a). On the other hand, methods that use 417 pattern information from the training models (methods 18-30; shown with diamond symbols in 418 the Taylor diagrams) estimate forced trends that look more like the ensemble-mean of the training 419 models (Fig. 1g-i, cf. Fig. 1d; 2g-i, cf. Fig. 2d). This is especially true for SST, and we suspect 420 that the reason for more diversity in forced precipitation trends is that not all training models have 421 the same forced precipitation response. Methods that use pattern information generally perform 422 better in terms of nRMSE than the methods that do not, but they will be more influenced by 423 any systematic biases in the training models, and they do not perform as well in terms of pattern 424 correlation for precipitation. 425

The Taylor diagrams for 1980-2022 trends in all 8 variables are shown in Fig. 3 and 4. For all 434 variables, the majority of ForceSMIP methods are skillful, where we consider a method skillful if 435 $\delta RMSE_i/RMSE_{RAW} < \delta r_i/r_{RAW}$, i.e., if the fractional reduction (improvement) in RMSE compared 436 to the raw data is greater than any fractional reduction (deterioration) in pattern correlation (below 437 the white lines in Fig. 3 and 4). Hence, a skillful method is required to reduce $RMSE_i$ compared 438 to $RMSE_{RAW}$, while at the same time not deteriorating the pattern correlation too strongly. This 439 definition of "skillfulness" thus implements the trade-off seen for some variables, such as precipi-440 tation, where a reduction in RMSE may be compensated by a deterioration of pattern correlation. 441 Skill for SST, T2m, TXx, and TNn are similar in an absolute sense, with $nRMSE_i$ between 0.3 442 and 0.6 (i.e., 30-60% errors). However, there is more improvement compared to the raw data 443 for TNn than for the other three surface-temperature variables, due to the larger signal-to-noise 444 ratio of TNn changes (not shown). The most skillful methods are generally similar across the 4 445 surface-temperature variables (i.e., methods 22, 23, 24, 25). There also tends to be a cluster of 446 simple methods with modest but systematic improvement compared to the raw data. 447



FIG. 3. Taylor diagram of method skill for 1980-2022 trends in (a) SST, (b) surface air temperature, (c) precipitation, and (d) sea level pressure. Colors, lines, and symbols as described in Fig. 1. Outlier methods excluded from the plots are: (a) 9, 30; (b) 27; (c) none; (d) none.

The absolute skill of the methods for trends in PR, SLP, and Rx1day is lower than for the four surface-temperature variables (Figs. 3c,d, 4c; cf. Figs. 3a,b, 4a,b). However, the improvement in nRMSE compared to the raw data is much larger for these variables. This occurs because there is a larger internal variability contribution to the 1980-2022 trends in these variables, and simply



FIG. 4. Taylor diagram of method skill for 1980-2022 trends in (a) annual maximum daily maximum temperature (TXx), (b) annual minimum daily minimum temperature (TNn), (c) annual maximum daily precipitation (Rx1day), and (d) zonal-mean atmospheric temperature (zmTa). black lines, and symbols as described in Fig. 1. Outlier methods excluded from the plots are: (a) 13, 27; (b) none; (c) none; (d) 9, 14, 15, 25. Note additionally that methods 1, 13, 16, 20, 21, 24, 27, and 30 did not estimate the forced response in zmTa.

reducing the amplitude of the raw data would reduced nRMSE. Some of the ML methods (e.g.,
25, 27) and one of the fingerprinting methods (24) even take the extreme approach of reducing

the estimated forced response amplitude to near zero for these variables, which does nevertheless 454 reduce nRMSE. The ability to improve nRMSE simply by reducing the amplitude of the estimated 455 forced trend pattern means that we should also pay attention to pattern correlation, which is 456 not influenced by the amplitude. Several of the simple methods consistently improve pattern 457 correlation across these variables (e.g., 6, 7, 8, 12, 16), as does one neural network method (20). 458 Of all variables, annual-mean precipitation (PR) shows the largest number of methods that reduce 459 the pattern correlation compared to the raw data, illustrating the difficulty in isolate the forced 460 response for this variable. 461

The skill for zonal-mean atmospheric temperature (zmTa) trends is an interesting case, because here the trend in the raw data is already such a skillful estimate of the forced response (nRMSE_{RAW} < 0.25) that only about half the methods can improve the skill further for this variable.

Here, we have focused on 1980-2022 trends, due in part to recent literature about SST trends over this time period (e.g., Wills et al. 2022). However, we also evaluated skill for other time periods, and the skill for 1950-2022 and 2000-2022 trends in SST are compared to the skill for 1980-2022 trends in Fig. S2. Methods generally show comparable absolute skill across the three time periods, however this represents a much larger improvement compared to the raw data for the short-term trends (2000-2022). This shows that the ForceSMIP methods have even more added value for short-term trends.

To more easily compare across methods and variables, Fig. 5 shows a scorecard for the two 480 main skill metrics, nRMSE_i and uncentered pattern correlation r_i . 1 – nRMSE_i is shown in place 481 of $nRMSE_i$ so that increased skill is positive in both panels. No single method stands out as most 482 skillful across all variables. While the fingerprinting and ML methods that use pattern information 483 from the training models (i.e., methods 18-30) generally stand out in terms of nRMSE, they tend 484 to have lower pattern correlation than simple methods (especially methods 1-8, 12, and 16). The 485 too low amplitude of some ML estimates is not apparent here, so it is important to keep in mind 486 the Taylor diagrams as well (cf. Figs. 3 and 4). There are a number of methods that have problems 487 with specific variables despite skill in other variables. One more general problem is the failure of 488 dynamical mode methods (e.g., 10, 11, 13, 14, 15) applied directly to variables such as PR, SLP, 489 and Rx1day that do not have the monthly or longer autocorrelation that generally underlies these 490 methods. An apparently successful workaround is to apply the dynamical mode method to SST or 491



FIG. 5. Skill summary scorecards for all methods' skill in 1980-2022 trends in all variables: (a) one minus 472 the normalized RMSE, normalized by the amplitude of the forced response, as in the Taylor diagrams; (b) the 473 uncentered pattern correlation. The root mean square nRMSE and average uncentered pattern correlation are 474 computed over the 5 "unseen model" evaluation members. Grey indicates that the method did not include a 475 forced response estimate for zmTa. Stippling indicates metrics where the ForceSMIP method gives a more 476 skillful forced trend estimate than the raw data, where the skill of estimating the forced trend by the raw data is 477 shown on the left hand side for reference. Note that values less than -1 in (a) are cropped and the colorbar in (b) 478 increases linearly with the square of the correlation. 479

⁴⁹² another variable with large autocorrelation and then to use the covariance with other variables to ⁴⁹³ get the forced response in the other variables, as was done by methods 8 and 12. Methods that stand ⁴⁹⁴ out in terms of consistency, with a consistent skill improvement relative to the raw data (stippling ⁴⁹⁵ in Fig. 5), are 7, 20, 24, 26, and 29, which includes the TrainingEM reference method (26). The ⁴⁹⁶ absolute skill of methods varies based on which evaluation member they are applied to (Fig. S1), ⁴⁹⁷ but the relative skill of methods relative to one another stays roughly the same across evaluation ⁴⁹⁸ members.

⁴⁹⁹ b. Spatiotemporal variability and large-scale climate indices

The long-term trends are only one way to evaluate the forced response estimates from the ForceSMIP methods, which include full spatiotemporal variability over 1950-2022. In this section we consider their skill for the spatiotemporal variability in the forced response, both at the grid scale and in selected large-scale climate indices.

⁵¹² We first synthesize the ForceSMIP methods' skill for grid-scale annual-mean spatiotemporal ⁵¹³ variability. Figure 6a shows 1 – nRMSE, where nRMSE is the global-mean RMSE in the grid-



FIG. 6. Skill summary scorecards for all methods' globally averaged skill in 10-yr running-mean grid-point 504 variability in all variables: (a) one minus the normalized RMSE, normalized by the amplitude of the forced 505 response; (b) square-root of global-mean correlation squared. The root mean square nRMSE and average 506 correlation are computed over the 5 "unseen model" evaluation members. Grey indicates that the method did 507 not include a forced response estimate for zmTa. Stippling indicates metrics where the ForceSMIP method has 508 more skill than the raw data, where the skill of estimating the forced response by the raw data is shown on the 509 left hand side for reference. Note that values less than -1 in (a) are cropped and the colorbar in (b) increases 510 linearly with the square of the correlation. 511

scale forced response estimate normalized by the global-mean root-mean-square amplitude of 514 the true forced response, estimated by the ensemble mean of the corresponding large ensemble. 515 Figure 6b shows the global-mean grid-point correlation of the forced response estimate and the 516 corresponding true forced response (ensemble mean). The absolute skill in both of these skill 517 metrics is less than the absolute skill in long-term trends (cf. Fig. 5), however, the skill added by 518 the ForceSMIP methods compared to the raw data is larger, and there is more widespread stippling, 519 indicating improvement relative to the raw data. All methods show consistent improvement relative 520 to the raw data across all variables in nRMSE, with a few exceptions in zmTa. Methods 1, 6-8, 12, 521 16, 21, 25, 29, and 30 additionally show improvement relative to the raw data across all variables 522 (except zmTa) in correlation. The skill of methods relative to one another is overall quite similar 523 for the spatiotemporal variability as for the long-term trends. 524

To evaluate the ForceSMIP methods' skill for large-scale climate indices, we choose 6 example 525 indices: (1) Annual-mean global-mean surface air temperature (GMST), (2) annual-mean Niño3.4 526 SST minus global-mean SST (GMSST), (3) the North Atlantic SST index (NASSTI) of the AMV, 527 i.e., annual-mean SST averaged over 0-60°N, 0-80°W minus the global mean, (4) Sahel monsoon 528 precipitation in MJJAS, averaged over 10-20°N, 20°W-10°E, (5) DJF Aleutian low SLP averaged 529 over 30-65°N, 160°E-140°W, and (6) TXx averaged over Continental Europe (land in 40-55°N, 530 0-40°E). A 10-yr running-mean is applied to indices 2-5 to filter out some of the high-frequency 531 noise, which would otherwise persist even in the ensemble mean of a large ensemble. 532

The skill of the ForceSMIP methods for these six large-scale indices is shown in Fig. 7. In 539 general, there are larger and more systematic nRMSE reductions compared to the raw data than 540 for the long-term trends in the corresponding variables (cf. Figs. 3 and 4). While there is 541 improvement in the correlation skill compared to the raw data for almost all methods in GMST 542 and Continental Europe TXx, there is more varied correlation skill across methods in the other 543 four indices. However, for each index, there is a subset of methods that are substantially improving 544 skill in terms of both nRMSE and correlation. Methods that consistently add skill compared to the 545 raw data across all indices (3-8, 12, 14-16, 18, 22, 24, 25, and 29) include a wide range of method 546 types, including both simple and complex methods. 547



FIG. 7. Taylor diagram showing skill for annual-mean temporal variability of climate indices: (a) GMST,
(b) 10-year running-mean Niño3.4 SST minus global mean SST, (c) 10-year running-mean NASSTI SST minus
global mean SST, (d) 10-year running-mean MJJAS Sahel precipitation, (e) 10-year running-mean DJF Aleutian
Low SLP, and (f) continental Europe (40-55°N, 0-40°W) TXx. Colors, lines, and symbols as described in Fig.
1, except with pattern nRMSE and pattern correlation replaced with nRMSE and correlation in these indices.
Outlier methods excluded from the plots are: (a) 13, 27, (b) 1, (c) none, (d) 20, 26, (e) none, (f) none.

548 5. Estimating the Forced Response in Observations

The underlying motivation for comparing StatML methods within ForceSMIP is to improve estimates of the forced response in observations. Now, armed with knowledge about which methods are skillful for which variables and metrics, we are ready to estimate the forced response in observations.

Each ForceSMIP method was applied to observational/reanalysis data in the same way it was 553 applied to the evaluation members used for method evaluation in the previous section. Our goal 554 in this section is to provide some examples of the observational forced responses estimated by 555 the ForceSMIP methods; a follow-up paper will use method weighting to generate a definitive 556 ForceSMIP forced response estimate with uncertainties. It is illustrative to first examine the forced 557 responses for individual skillful methods. In Figs. 8, 9, and 10, we show the forced and internal 558 components of observed 1980-2022 trends in SST, PR, and SLP, respectively, as estimated by 559 selected ForceSMIP methods, alongside the raw observed trends over this period. Methods are 560 selected to illustrate the range of different forced trend estimates, based on an EOF analysis in 561 Appendix A. 562

The strong pattern observed in the 1980-2022 SST trend, with cooling in the East Pacific and 565 Southern Ocean and intensified warming in the West Pacific and North Atlantic, unlike the more 566 uniform East-Pacific intensified warming the climate models show for this period, has generated 567 substantial interest from the climate science community (Wills et al. 2022; Seager et al. 2022; 568 Watanabe et al. 2024; Simpson et al. 2025). This lack of agreement with models is apparent in 569 the comparison in Fig. 8 with the TrainingEM method (26), which is equal (up to an amplitude 570 rescaling) to the ensemble mean of the 5 training models. The residual internal variability estimated 571 by TrainingEM is large, and has been shown to be larger than is consistent with internal variability 572 in most climate models (Wills et al. 2022; Seager et al. 2022). 573

Several of the other ForceSMIP methods shown have a smaller amplitude of estimated internal variability in 1980-2022 SST trends, indicating that they are estimating a forced response that is closer to the full observed trends than is the TrainingEM forced response. However, the degree to which individual methods' forced response estimates are more similar to the full observed trends or to the TrainingEM forced response varies substantially. LFCA-2 is one end member, estimating that almost all of the observed trend over 1980-2022 is forced. EOF-SLR is another end member, with Observed SST Trend (1980-2022) (°C per 42 yr)



FIG. 8. Forced and internal components of observed SST trends (1980-2022) for TrainingEM and selected skillful methods, chosen as representative examples from the EOF analysis in Figure A1.

⁵⁸⁰ a forced response similar to TrainingEM except for reduced El-Niño-like warming and somewhat ⁵⁸¹ more warming in the Atlantic. GPCA and UNet3D-LOCEAN are in between these end members, ⁵⁸² but each with their own unique features. The differences across these methods, all of which ⁵⁸³ are shown to be skillful in the method evaluation (Fig. 3a), illustrates the epistemic uncertainty ⁵⁸⁴ in estimating the forced response from observations, where epistemic uncertainty refers to the ⁵⁸⁵ uncertainty and potential systematic biases associated with the method used for forced response ⁵⁸⁶ estimation. While EOF-SLR and UNet3D-LOCEAN are modestly more skillful than the other ⁵⁸⁷ methods in the method evaluation, we cannot say with certainty which of these six forced response ⁵⁸⁸ estimates is closer to the truth.



FIG. 9. Forced and internal components of observed PR trends (1980-2022) for TrainingEM and selected skillful methods, chosen as representative examples from the EOF analysis in Figure A2.

There is even wider spread of forced response estimates for precipitation (Fig. 9; see also Fig. 591 A2), ranging from MF-LFCA-2 estimating that most of the observed 1980-2022 trend is forced to 592 MonthFinger and TrainingEM estimating that almost none of it is. MF-LFCA and SNMP-OF are 593 somewhere in between, with forced and internal contributions of similar amplitudes. It is worth 594 noting that by focusing on forced responses that are robust across models, the estimated forced 595 responses by TrainingEM and MonthFinger are smaller in amplitude than the forced precipitation 596 response in individual models (cf. Fig. 2b), due to structural differences in models' forced 597 responses. 598

The estimated 1980-2022 forced trends in SLP are all quite different from one another (Fig. 601 10). They agree on the poleward shift of the Southern Hemisphere westerly winds indicated by 602 the positive and negative bands of SLP trends north and south of ~ 50°S, but they have more 603 than a factor of four spread in the magnitude of this circulation change. Some methods show 604 that the Aleutian low weakening is mostly forced (MF-LFCA-2, consistent with the SST estimate 605 from LFCA-2 in Fig. 8) while others show it is almost entirely internal variability (MF-LFCA, 606 UNet3D-LOCEAN, ANN-Fingerprinters). There is a similar lack of agreement on whether North 607 Atlantic SLP trends are forced or unforced. The large uncertainty in the forced response of SLP is 608 consistent with the literature (Knutson and Ploshay 2021). 609

To get a sense for the average separation of 1980-2022 trends into forced and internal components 610 by the ForceSMIP methods, we average the forced response estimates over all ForceSMIP methods 611 determined to be skillful for each variable. Methods are included if the improvement in RMSE 612 exceeds the deterioration of pattern correlation ($\delta RMSE_i/RMSE_{RAW} < \delta r_i/r_{RAW}$; below the white 613 lines in Fig. 3 and 4). Figs. 11 and and S3 show the forced trend averaged over these skillful 614 methods, the residual internal variability component of the trends, and the forced trend estimated 615 by TrainingEM, which gives a sense of what climate models say the forced response should be 616 over this time period. 617

The ForceSMIP-skillful-method-mean (hereafter ForceSMIP-mean) forced SST trend over 1980-2022 shows near-zero warming in the East Pacific and South Pacific, where the full observed SST trend shows cooling. The ForceSMIP-mean therefore attributes some but not all of the difference in 1980-2022 SST trend pattern between models and observations to internal variability. Similarly, the observed cooling of the Southern Ocean, which is not reproduced by models, is attributed Observed SLP Trend (1980-2022) (Pa per 42 yr)



FIG. 10. Forced and internal components of observed SLP trends (1980-2022) for TrainingEM and selected skillful methods, chosen as representative examples from the EOF analysis in Figure A3.

to a combination of forced response and internal variability. The ForceSMIP-mean also shows stronger weakening of the Aleutian Low and stronger strengthening of the Amundsen Sea Low than TrainingEM, which are both similar to La Niña teleconnections. ForceSMIP also suggests a more La–Niña-like forced trend in precipitation, with a much larger amplitude than the estimate by TrainingEM. However, as noted previously, the TrainingEM estimate for precipitation is smaller



FIG. 11. (center column) Mean estimates of forced trends (1980-2022) over all skillful ForceSMIP methods (defined as δ RMSE_{*i*}/RMSE_{RAW} < $\delta r_i/r_{RAW}$, i.e., below the white line in Figs. 3 and 4) for SST, SLP, precipitation, TXx, TNn, and Rx1day. Units are °C per 42 yr, Pa per 42 yr, or mm day⁻¹ per 42 yr accordingly. (right column) The residual trends attributed to internal variability. (left column) The TrainingEM reference method, obtained from the multi-model-mean of the five training models, is shown for comparison.

than the forced response in individual models because it focuses on the common response across
all 5 training models.

The ForceSMIP-mean 1980-2022 forced trends in T2m, TXx, and TNn are broadly similar over 635 ocean regions (Figs. 11 and S3), where they show a more La-Niña-like forced response than 636 TrainingEM and less warming in the Kuroshio-Oyashio extension, consistent with what was found 637 for SST. The forced trend in TXx shows more warming than the forced trend in T2m in tropical 638 land regions and less in high-latitude land regions, whereas the opposite is true for the forced trend 639 in TNn. This is consistent with the reduction (increase) in temperature variability in high-latitude 640 (tropical) land regions (Kotz et al. 2021), and is also seen in TrainingEM. TXx and TNn both 641 have larger estimated contributions of internal variability to 1980-2022 trends than does T2m, 642 illustrating the added value of the ForceSMIP methods for noisy extreme-event statistics. Rx1day 643 has by far the largest estimated contribution of internal variability to 1980-2022 trends, though the 644 estimated forced response is still larger than that estimated from TrainingEM. 645

To visualize the ForceSMIP-estimated forced responses in the six climate indices, Figure 12 shows the likely (66%) range (i.e., the 17th and 83rd percentiles) of the ForceSMIP methods determined to be skillful compared to the raw observed data , as well as TrainingEM and five example methods. Methods are considered skillful and thus included in the likely range if they show a fractional reduction in nRMSE that exceeds any fractional reduction in their correlation (below the white lines in Fig. 7). Example methods are chosen that have varying complexity, high skill across most variables, and produce different forced response estimates from one another.

⁶⁵⁷ Compared to the raw data, all skillful methods smooth out some of the interannual variability in ⁶⁵⁸ GMST (Fig. 12a). On a quantitative level, there is a 66% uncertainty range in the estimated forced ⁶⁵⁹ 1950-2022 GMST trend of 0.89-1.07°C per 72 yr. The smoothing out of interannual variability ⁶⁶⁰ is even more important for metrics such as Continental Europe TXx, where the forced responses ⁶⁶¹ estimates are all much smoother than the raw data (Fig. 12f). Methods consistently attribute the ⁶⁶² multi-year negative excursion between 1975 and 1980 to internal variability. The ratio of estimated ⁶⁶³ forced trends in Continental Europe TXx and GMST has a 66% range of 1.89-2.79.

⁶⁶⁴ While the forced responses in GMST and Continental Europe TXx could be guessed to some ⁶⁶⁵ degree of accuracy by simply smoothing the raw data, estimating the forced components of the ⁶⁶⁶ other four indices is much more challenging. The ForceSMIP estimated observed forced response in 10-yr running-mean Niño3.4 (minus GMSST) ranges from increasing (El-Niño-like warming)
in Anchor-OPLS and TrainingEM to monotonically decreasing (La-Niña-like warming) in MFLFCA and SNMP-OF (Fig. 12b), with MF-LFCA-2 even showing a strong increase through 1980
followed by a strong decrease. Nevertheless, all methods agree that the large negative excursion in
the early 1970s and the large positive excursion in the early 1990s resulted from internal variability.
The 66% range in the estimated 1950-2022 forced trend in Niño3.4 minus GMSST is -0.27-0.10°C
per 72 yr, indicating that even the sign of the long-term forced trend remains uncertain.

The estimates of how much the AMV is forced range from almost all of it to none of it, as well as everything in between (Fig. 12c). ForceSMIP thus helps to reconcile research that indicates



FIG. 12. Raw data timeseries, scaled training models ensemble mean (TrainingEM), skillful methods likely (66%) range, and selected methods for each index in Figure 7. Skillful methods are defined as those with a fractional reduction in nRMSE that exceeds any fractional reduction in their correlation (below the white lines in Fig. 7).

that the AMV is mostly forced (Booth et al. 2012; Wills et al. 2020; He et al. 2023) with research 676 suggesting that it is mostly internal variability (Ting et al. 2009; Zhang et al. 2013; Qin et al. 677 2020; Latif et al. 2022) by demonstrating that either could be true. Interestingly, the two end 678 members with most and least forced AMV are MF-LFCA and MF-LFCA-2, which differ only in 679 the number of low-frequency patterns included, this illustrates how the hyperparameter sensitivity 680 of the LFCA method may actually help to quantify the epistemic uncertainty in the forced response 681 estimate. Given the association between the AMV and Sahel precipitation (Zhang and Delworth 682 2006), it is not surprising that there is also a large spread in the forced response estimates for 683 Sahel precipitation (Fig. 12d). What is interesting however is that all of the ForceSMIP estimates 684 either show a drying or a much weaker wettening trend than TrainingEM. This suggests that CMIP6 685 models, at least those used for training, have systematic biases in Sahel precipitation trends. Finally, 686 the ForceSMIP methods consistently show a small forced response in the Aleutian Low, attributing 687 its large decadal excursions to internal variability (Fig. 12e). 688

Overall, ForceSMIP provides an ensemble of estimates of the observed forced response, and we highlight cases where there are consistent differences from the forced response in climate models (e.g., the La-Niña-like forced response in observations) as well as cases where epistemic uncertainty limits the ability to draw conclusions (e.g., on the amplitude of forced AMV).

6. Conclusions, Discussion, and Outlook

We have demonstrated that many different types of StatML methods exhibit skill in estimating 694 the forced response from individual ensemble members of a climate model large ensemble, where 695 skill means that they give a better forced response estimate than the raw data. Skillful methods 696 include simple regression approaches, LFCA, LIM-based methods, as well as fingerprinting and 697 ML methods custom built for the ForceSMIP project. Methods are most skillful in absolute terms 698 for thermodynamic responses, such as in SST and surface air temperature, but the added value 699 of these methods compared to the raw data is largest for responses in fields with large amplitude 700 internal variability such as SLP, precipitation, and extreme-event indices. The ForceSMIP methods 701 are skillful for long-term regional-scale trends (e.g., over 1980-2022), grid-scale spatiotemporal 702 variability, and large-scale climate indices. No single method outperforms the others across all 703 variables, but rather the most skillful methods vary depending on the metric of evaluation. 704

Armed with an array of skillful methods for forced response estimation, we investigated the 705 forced response in observations in Section 5. We found that the ForceSMIP methods systematically 706 estimate that the observed forced response is more La-Niña-like than indicated by models, with a 707 local minimum in warming in the Southeast Pacific, but also that the discrepancy in 1980-2022 SST 708 trends between observations and models is partly due to internal variability. The observed forced 709 response obtained from the average of skillful ForceSMIP methods also exhibits La-Niña-like 710 teleconnections in other variables, including SLP and precipitation. Despite these commonalities, 711 there is a large spread in the estimated forced SST trend pattern across methods that display similar 712 skill in the large ensemble evaluation data, and an even wider spread of forced responses for SLP 713 and precipitation. The spread across estimates of the forced response is sufficiently large that most 714 statements about the relative contributions of external forcing and internal variability (for example 715 to the AMV) cannot be made with great certainty. Overall, ForceSMIP suggests that there are 716 systematic differences in the forced response between climate models and observations while also 717 illustrating the intrinsic epistemic uncertainty in estimating the forced response from observations. 718 The intrinsic uncertainty in the extent to which multi-decadal SST fluctuations and regional details 719 of trend patterns are forced or unforced is important to consider in the context of climate change 720 attribution, model evaluation, and climate impact assessments. 721

722 a. Which method should I use?

At this point, you may be wondering, which method should I use for forced response estimation in my own work? While the method evaluation in Figs. 3-7 may give some guidance, it's quite likely that this paper did not consider your metric of interest. Furthermore, there are almost always many good choices for any given metric. Nevertheless, we can give a few recommendations:

Use more than one type of method to get a better sense of how the forced response estimate
 varies across methods. It's worth keeping in mind that simple methods tend to stay closer
 to the observed trends, whereas most fingerprinting and ML methods will give observational
 forced response estimates more similar to the forced response in the climate models used for
 training, and will thus be more subject to any systematic biases in the training dataset.

⁷³² 2. Either use methods that generalize well across metrics or train/test the methods you use
 ⁷³³ for your metric of interest within a large ensemble dataset. The diversity of variables and

metrics considered by ForceSMIP makes it likely that methods consistently showing skill in
 ForceSMIP will generalize well to other applications.

736 737 3. The ForceSMIP evaluation dataset (Wills et al. 2025) is a useful resource for evaluating new methods and/or for evaluating which methods work best for a specific application of interest.

Finally, another relevant consideration is that the ML methods would all need to be re-trained for
other applications, whereas most of the other methods work out of the box and do not need further
customization. However, the need to train ML methods can also be an advantage, because it means
they will be tailored for the application of interest.

742 b. Lessons for further method development

Several lessons can be learned from the successes and failures of individual ForceSMIP methods. 743 One of the clearest lessons is that - perhaps to no great surprise - LIMs only perform well 744 for variables that have sufficiently large autocorrelation on the timescale of interest (monthly 745 anomalies in our case). This is exemplified by the much higher skill of LIMnMCA and Colored-746 LIMnMCA compared to other LIM-based methods for variables such as precipitation, SLP, and 747 RX1day. What's different about these two methods is that they applied a LIM to SST and then 748 used maximum covariance analysis to identifying the covarying forced patterns in other variables. 749 Another approach could be to merge each field variable with SST and apply a joint analysis to 750 both fields at once. This approach was used for MF-LFCA, where it led to modest improvement 751 in skill for precipitation and SLP over the one-field-at-time LFCA. We highlight these cases due 752 to the clean comparisons they offer, but several other methods used multiple fields at once (Table 753 1). Many of the methods that analyzed one field variable at a time could likely be improved by 754 applying them to two or more field variables at a time, especially if the additional variable is a field 755 with a clear forced response, such as SST. 756

Another lesson is that methods focused on reducing RMSE or related metrics may end up guessing a near-zero forced response in cases where internal variability is larger than the forced response. To control against this, methods could expand the skill metrics they consider, for example by incorporating correlation or amplitude-error metrics and computing skill metrics on different timescales. This could draw on the experiences of the machine-learning weather prediction community (e.g., Nathaniel et al. 2024), which is grappling with similar issues. Some methods may also give better forced response estimates if they were reformulated to explicitly estimate
 both forced and unforced climate variations, as was already done in UNet3D-LOCEAN (see also
 Po-Chedley et al. 2022).

An additional important consideration for further method development is that the ML methods are by design more trainable to optimize for a specific task. We intentionally did not specify exact evaluation targets in advance for this phase of ForceSMIP, to avoid all methods overfitting to particular metrics. Further development of these methods can now focus on correcting for some of the problems displayed in this round of evaluation. Future work should focus on cataloging a comprehensive set of forced response metrics of interest, so that methods can be trained to optimize across many relevant metrics at once.

c. An observational forced response estimate and its applications

A primary goal of ForceSMIP is to generate a forced response in observations, including a 774 quantification of the associated epistemic uncertainty, i.e., uncertainty from different methods 775 of estimation getting different answers. In this study, we have provided one such estimate: a 776 30-method ensemble of different forced response estimates (openly available on Zenodo; Wills 777 et al. 2025). We additionally quantified the expected error based on evaluation within large 778 ensembles and gave demonstrations of the types of information that can be obtained from such a 779 multi-method ensemble, showing both differences in the estimated forced response across methods 780 as well as the multi-method-mean forced response estimate for skillful methods. The method 781 weighting is intentionally kept simple in this paper, with methods given full weight for skill above a 782 threshold and zero weight otherwise. A follow-up paper will apply a systematic method weighting 783 scheme, following Merrifield et al. (2023), to provide a skill weighted forced response estimate 784 and uncertainty range. We also encourage others to generate their own forced response estimates 785 from this dataset that are customized to specific applications. 786

We foresee many possible applications of an observational forced response estimate. One set of applications is for model evaluation. An observational forced response from ForceSMIP could be combined with an estimate of the residual variance due to estimation uncertainty and internal variability, e.g., based on the nRMSE evaluated in Section 4, and this would then provide a comparison point for evaluating forced trends in models against observations (cf. Simpson et al.

2025). The flip-side of evaluating forced trends in models is evaluating their amplitude of internal 792 decadal variability, which has been suggested based on instrumental and paleoclimate data to be 793 too weak in some regions (Laepple and Huybers 2014; Dee et al. 2017; Laepple et al. 2023). 794 ForceSMIP can help to evaluate whether there are discrepancies in forced or internal multi-decadal 795 variance compared to large ensembles. However, our results already suggest that, for metrics 796 with large multi-decadal variability such as the AMV, the separation between forced and internal 797 components remains extremely challenging, with some methods estimating a forced response more 798 like the raw observations and some methods estimating a forced response more like the ensemble 799 mean of the training models. In these cases, it will remain difficult to distinguish between model 800 discrepancies in the forced response and model discrepancies internal variability. 801

Another set of applications of forced response estimates from ForceSMIP is for monitoring 802 internal climate variability and generating observational large ensembles. Indices of internal 803 variability, where the forced response is often removed either by removing the linear trend or by 804 subtracting GMSST, can increasingly be influenced by climate change. For example, Deser and 805 Phillips (2023b) show how not fully removing the forced response from indices of the AMV can lead 806 to spurious implied connections with the tropical Pacific. We therefore suggest that the ForceSMIP 807 forced response, if continuously updated, could serve as a standard estimate of the forced response 808 to remove from indices of internal variability such as ENSO, AMV, PDO, and NAO. Removal of 809 the forced response also allows for generation of an observational large ensemble, e.g., using the 810 phase randomization approach of McKinnon and Deser (2018, 2021). Such an observational large 811 ensemble can help to explore long-term trends and extreme events that could have happened in the 812 real world under different phasing of internal variability (e.g., as in Deser and Phillips 2023a). 813

Underlying all of these applications of ForceSMIP observational forced response estimates is the intrinsic interest in the observational forced response itself, which can help to understand and communicate how anthropogenic activities have affected historical climate and give a glimpse into the changes expected in the near future.

Acknowledgments. This research benefited greatly from synchronous in-person hackathons in
 Boulder, CO and Zurich, Swizerland in August 2023, which were funded by the U.S. National
 Science Foundation, the Swiss National Science Foundation (Award IZSEZ0-220740), the Inter national CLIVAR Project Office, and the Packard Foundation. R. C. J. Wills was supported by

the Swiss National Science Foundation (Award PCEFP2-203376). C. Deser and A. Phillips were 822 supported by the NSF National Center for Atmospheric Research, which is a major facility spon-823 sored by the NSF under the Cooperative Agreement 1852977. K. A. McKinnon was supported by 824 the Packard Foundation. S. Po-Chedley, C. Bonfils, S. Duan, and M. A. Fernandez were funded 825 by the Regional and Global Model Analysis program area of the U.S. Department of Energy's 826 (DOE) Office of Biological and Environmental Research (BER) as part of PCMDI, an Earth Sys-827 tem Model Evaluation Project. Work by S. Po-Chedley, C. Bonfils, and S. Duan was performed 828 under the auspices of the U.S. Department of Energy by Lawrence Livermore National Laboratory 829 under Contract DE-AC52-07NA27344. S. Sippel acknowledges the climXtreme project funded 830 by the German Federal Ministry of Education and Research (Phase 2, project PATTETA, Grant 831 No. 01LP2323C) and the EU Horizon project AI4PEX (Grant agreement No. 101137682). C. 832 Bône and G. Gastineau acknowledge the support of the EUR IPSL Climate Graduate School 833 project managed by the ANR under the "Investissements d'avenir" programme with the reference 834 ANR-11-IDEX-0004-17-EURE-0006. G. Camps-Valls, H. Durand, and G. Varando acknowledge 835 funding from the European Research Council (ERC) under the ERC Synergy Grant USMILE (grant 836 agreement 855187) and funding from the Horizon project AI4PEX (grant agreement 101137682). 837 N. Mankovich acknowledges support from the project "Artificial Intelligence for complex systems: 838 Brain, Earth, Climate, Society" funded by the Department of Innovation, Universities, Science, and 839 Digital Society, code: CIPROM/2021/56. J.-R. Shi was supported by U.S. National Science Foun-840 dation under Grant OCE-2048336. The EOF-SLR and LDM-SLR methods were developed under 841 the support of the state assignment of the Institute of Applied Physics of the Russian Academy 842 of Sciences (Project FFUF-2022-0008). The results from UNet3D-LOCEAN were performed 843 using HPC resources from GENCI-IDRIS AD011013295R2 and AD011013295R3. We would 844 like to acknowledge computing support from the Casper system (https://ncar.pub/casper) provided 845 by the NSF National Center for Atmospheric Research (NCAR), sponsored by the National Sci-846 ence Foundation. The authors thank all participants in the ForceSMIP hackathons for valuable 847 discussions. 848

Data availability statement. The CMIP6 source data are available via ESGF, and the processed
 large ensemble data used in ForceSMIP has recently been made available by Maher et al. (2025).
 ERA5 data is available from https://cds.climate.copernicus.eu/datasets. ERSSTv5 data is available

from https://psl.noaa.gov/data/gridded/data.noaa.ersst.v5.html. The ForceSMIP Tier 1 data, i.e., 852 the raw data, ensemble means, and estimated forced responses for each variables and each evalua-853 tion member is available on Zenodo (Wills et al. 2025). The code for all StatML methods is made 854 available via Github (https://github.com/ForceSMIP/tier1-methods). Scripts for evalu-855 ating methods using the ForceSMIP Tier 1 data are made available in a separate Github repository 856 (https://github.com/ForceSMIP/tier1-evaluation), including an example script for eval-857 uating forced trends in Python and all MATLAB scripts used for analysis in this paper. Permanent 858 identifiers for these Github repositories will be included in the final published version. 859

APPENDIX

860

861

Analysis of Inter-Method Variance

In order to illustrate the inter-method differences (i.e., epistemic uncertainty) in estimated forced 867 trends, we perform an EOF analysis on the forced trends estimated by skillful methods. Methods 868 are included if $\delta RMSE_i/RMSE_{RAW} < \delta r_i/r_{RAW}$ (below the white lines in Fig. 3). The results 869 are shown for the EOF analysis of estimated 1980-2022 forced trends in SST, PR, and SLP in 870 Figs. A1, A2, and A3, respectively. Panels (a) and (b) show the EOF patterns and the percentage 871 of the variance they explain. Panels (c) show the corresponding principal components, i.e., the 872 contribution of each EOF to the forced trend estimated by each method. The distribution of 873 principal components are used to inform the selection of methods shown in Figs. 8-10, which are 874 highlighted with red symbols in panels (c) of Figs. A1-A3. 875

Estimated 1980-2022 forced trends in SST differ from one another in a pattern (EOF1) similar to what has been called the Interdecadal Pacific Oscillation (IPO; Power et al. 1999), indicating that some methods estimate the IPO to be mostly forced, while others do not. Methods also differ in their estimates of the amount of forced warming in the Northern Hemisphere ocean basins (EOF2). The net result is that there is uncertainty in the forced SST trend in the East Pacific, Southern Ocean, Kuroshio-Oyashio Extension, and subpolar North Atlantic (Fig. A1d).

The EOF analysis for estimated 1980-2022 forced trends in PR (Fig. A2) shows a large fraction of variance explained by EOF1, which resembles the full observed trend (Fig. 9). The amplitude of PC1 shows clusters near -1 and 1.5 (Fig. A2c), which are methods estimating that very little or most of the observed trend is forced, respectively. The leading EOF of estimated 1980-2022 forced trends in SLP (Fig. A3a) includes positive anomalies in the Aleutian low region and South Pacific and negative anomalies around Antarctic, resembling the SLP pattern associated with the IPO. Combined with EOF2 (Fig. A3b), the net result is uncertainty in the midlatitudes in all ocean basins as well as around Antarctica (Fig. A3d).

References

- Alexander, M. A., L. Matrosova, C. Penland, J. D. Scott, and P. Chang, 2008: Forecasting Pacific
- ⁸⁹⁶ SSTs: Linear inverse model predictions of the PDO. *Journal of Climate*, **21** (**2**), 385–402.
- Bellucci, A., A. Mariotti, and S. Gualdi, 2017: The role of forcings in the twentieth-century North
- Atlantic multidecadal variability: The 1940–75 North Atlantic cooling case study. Journal of
- ⁸⁹⁹ *Climate*, **30** (**18**), 7317–7337.



FIG. A1. Inter-method EOF analysis of estimated forced SST trends over 1980-2022, including only skillful methods (defined as δ RMSE/RMSE_{RAW} < $\delta r/r_{RAW}$, i.e., below the white line in Figs. 3 and 4). (a) Inter-method EOF1, (b) inter-method EOF2, and (c) the principal component amplitudes for each method. The percentage of total variance explained by each EOF is shown in the title of (a) and (b). (d) Total inter-method variance, expressed as a standard deviation. Red symbols in (c) indicate methods shown in Fig. 8.

- Bethke, I., and Coauthors, 2021: NorCPM1 and its contribution to CMIP6 DCPP. *Geoscientific Model Development*, 14 (11), 7073–7116.
- ⁹⁰² Blackport, R., and J. C. Fyfe, 2022: Climate models fail to capture strengthening wintertime North ⁹⁰³ Atlantic jet and impacts on Europe. *Science Advances*, **8** (**45**), eabn3112.
- Bône, C., G. Gastineau, S. Thiria, P. Gallinari, and C. Mejia, 2024: Separation of internal and
 forced variability of climate using a U-Net. *Journal of Advances in Modeling Earth Systems*,
 16 (6), e2023MS003 964.
- Booth, B. B., N. J. Dunstone, P. R. Halloran, T. Andrews, and N. Bellouin, 2012: Aerosols
 implicated as a prime driver of twentieth-century North Atlantic climate variability. *Nature*,
 484 (7393), 228–232.
- ⁹¹⁰ Boucher, O., and Coauthors, 2020: Presentation and Evaluation of the IPSL-CM6A-LR Cli⁹¹¹ mate Model. *Journal of Advances in Modeling Earth Systems*, **12** (7), e2019MS002010,
 ⁹¹² https://doi.org/doi.org/10.1029/2019MS002010.



FIG. A2. Same as A1, but for estimated forced PR trends over 1980-2022. Red symbols in (c) indicate methods shown in Fig. 9.



FIG. A3. Same as A1, but for estimated forced SLP trends over 1980-2022. Red symbols in (c) indicate methods shown in Fig. 10.

Brunner, L., M. Hauser, R. Lorenz, and U. Beyerle, 2020: The ETH Zurich CMIP6 next generation
archive: Technical documentation. Zenodo, https://doi.org/10.5281/zenodo.3734127.

⁹¹⁵ Compo, G. P., and P. D. Sardeshmukh, 2010: Removing ENSO-related variations from the climate
 ⁹¹⁶ record. *Journal of Climate*, 23 (8), 1957–1978.

⁹¹⁷ Dai, A., J. C. Fyfe, S.-P. Xie, and X. Dai, 2015: Decadal modulation of global surface temperature
⁹¹⁸ by internal climate variability. *Nature Climate Change*, 5 (6), 555–559.

Dee, S., L. Parsons, G. Loope, J. Overpeck, T. Ault, and J. Emile-Geay, 2017: Improved spectral
 comparisons of paleoclimate models and observations via proxy system modeling: Implications
 for multi-decadal variability. *Earth and Planetary Science Letters*, **476**, 34–46.

DelSole, T., M. K. Tippett, and J. Shukla, 2011: A significant component of unforced multidecadal
 variability in the recent acceleration of global warming. *Journal of Climate*, 24 (3), 909–926.

- Delworth, T. L., and Coauthors, 2020: SPEAR: The next generation GFDL modeling system for
 seasonal to multidecadal prediction and projection. *Journal of Advances in Modeling Earth Systems*, 12 (3), e2019MS001895, https://doi.org/doi.org/10.1029/2019MS001895.
- Deser, C., A. Phillips, V. Bourdette, and H. Teng, 2012: Uncertainty in climate change projections:
 the role of internal variability. *Climate dynamics*, 38, 527–546.
- Deser, C., and A. S. Phillips, 2021: Defining the internal component of Atlantic multidecadal
 variability in a changing climate. *Geophysical Research Letters*, 48 (22), e2021GL095 023.
- Deser, C., and A. S. Phillips, 2023a: A range of outcomes: the combined effects of internal variability and anthropogenic forcing on regional climate trends over Europe. *Nonlinear Processes in Geophysics*, **30** (1), 63–84.
- ⁹³⁴ Deser, C., and A. S. Phillips, 2023b: Spurious Indo-Pacific connections to internal Atlantic
 ⁹³⁵ Multidecadal variability introduced by the global temperature residual method. *Geophysical* ⁹³⁶ *Research Letters*, **50** (3), e2022GL100 574.
- ⁹³⁷ Deser, C., A. S. Phillips, M. A. Alexander, and B. V. Smoliak, 2014: Projecting North American
 ⁹³⁸ climate over the next 50 years: Uncertainty due to internal variability. *Journal of Climate*, 27 (6),
 ⁹³⁹ 2271–2296.
- Deser, C., and Coauthors, 2020: Insights from Earth system model initial-condition large ensembles
 and future prospects. *Nature Climate Change*, **10** (**4**), 277–286.
- Eyring, V., S. Bony, G. A. Meehl, C. A. Senior, B. Stevens, R. J. Stouffer, and K. E. Taylor, 2016:
 Overview of the Coupled Model Intercomparison Project Phase 6 (CMIP6) experimental design
 and organization. *Geoscientific Model Development*, 9 (5), 1937–1958.
- Frankcombe, L. M., M. H. England, M. E. Mann, and B. A. Steinman, 2015: Separating internal
 variability from the externally forced climate response. *Journal of Climate*, 28 (20), 8184–8202.
- Frankignoul, C., G. Gastineau, and Y.-O. Kwon, 2017: Estimation of the SST response to anthropogenic and external forcing and its impact on the Atlantic multidecadal oscillation and the
 Pacific decadal oscillation. *J. Climate*, **30** (24), 9871–9895.

- ⁹⁵⁰ Gavrilov, A., S. Kravtsov, M. Buyanova, D. Mukhin, E. Loskutov, and A. Feigin, 2024: Forced
 ⁹⁵¹ response and internal variability in ensembles of climate simulations: Identification and analysis
 ⁹⁵² using linear dynamical mode decomposition. *Climate Dynamics*, **62** (3), 1783–1810.
- ⁹⁵³ Gavrilov, A., S. Kravtsov, and D. Mukhin, 2020: Analysis of 20th century surface air temperature
 ⁹⁵⁴ using linear dynamical modes. *Chaos: An Interdisciplinary Journal of Nonlinear Science*,
 ⁹⁵⁵ **30** (12).
- Hajima, T., and Coauthors, 2020: Development of the MIROC-ES2L Earth system model and
 the evaluation of biogeochemical processes and feedbacks. *Geoscientific Model Development*, **13 (5)**, 2197–2244, https://doi.org/10.5194/gmd-13-2197-2020.
- Hasselmann, K., 1979: On the signal-to-noise problem in atmospheric response studies.
- Hawkins, E., and R. Sutton, 2009: The potential to narrow uncertainty in regional climate predictions. *Bulletin of the American Meteorological Society*, **90** (8), 1095–1108.
- He, C., A. C. Clement, S. M. Kramer, M. A. Cane, J. M. Klavans, T. M. Fenske, and L. N.
 Murphy, 2023: Tropical Atlantic multidecadal variability is dominated by external forcing.
 Nature, 622 (7983), 521–527.
- Hegerl, G. C., H. von Storch, K. Hasselmann, B. D. Santer, U. Cubasch, and P. D. Jones, 1996:
 Detecting greenhouse-gas-induced climate change with an optimal fingerprint method. *Journal of Climate*, 9 (10), 2281–2306.
- Hersbach, H., and Coauthors, 2020: The ERA5 global reanalysis. Q. J. R. Meteorol. Soc., 146 (730),
 1999–2049.
- ⁹⁷⁰ Huang, B., and Coauthors, 2017: Extended Reconstructed Sea Surface Temperature, version 5
 (ERSSTv5): upgrades, validations, and intercomparisons. *Journal of Climate*, **30** (**20**), 8179–
 ⁹⁷² 8205.
- ⁹⁷³ Hyvärinen, A., and E. Oja, 2000: Independent component analysis: Algorithms and applications.
 ⁹⁷⁴ *Neural networks*, **13 (4-5)**, 411–430.

- ⁹⁷⁵ Knutson, T. R., and J. Ploshay, 2021: Sea level pressure trends: Model-based assessment of
 ⁹⁷⁶ detection, attribution, and consistency with CMIP5 historical simulations. *Journal of Climate*,
 ⁹⁷⁷ **34** (1), 327–346.
- ⁹⁷⁸ Kotz, M., L. Wenz, and A. Levermann, 2021: Footprint of greenhouse forcing in daily temperature
 ⁹⁷⁹ variability. *Proceedings of the National Academy of Sciences*, **118 (32)**, e2103294 118.
- ⁹⁸⁰ Kravtsov, S., C. Grimm, and S. Gu, 2018: Global-scale multidecadal variability missing in state-⁹⁸¹ of-the-art climate models. *npj Climate and Atmospheric Science*, **1** (**1**), 34.
- Laepple, T., and P. Huybers, 2014: Ocean surface temperature variability: Large model-data
 differences at decadal and longer periods. *Proceedings of the National Academy of Sciences*,
 111 (47), 16682–16687.
- Laepple, T., and Coauthors, 2023: Regional but not global temperature variability underestimated
 by climate models at supradecadal timescales. *Nature Geoscience*, 16 (11), 958–966.
- Latif, M., J. Sun, M. Visbeck, and M. Hadi Bordbar, 2022: Natural variability has dominated
 Atlantic meridional overturning circulation since 1900. *Nature Climate Change*, **12** (**5**), 455–460.
- Lehner, F., C. Deser, N. Maher, J. Marotzke, E. M. Fischer, L. Brunner, R. Knutti, and E. Hawkins,
 2020: Partitioning climate projection uncertainty with multiple large ensembles and CMIP5/6.
 Earth System Dynamics, **11** (2), 491–508.
- Lien, J., Y.-N. Kuo, H. Ando, and S. Kido, 2025: Colored linear inverse model: A data-driven
 method for studying dynamical systems with temporally correlated stochasticity. *Physical Review Research*, 7 (2), 023 042.
- Maher, N., and Coauthors, 2025: The updated Multi-Model Large Ensemble Archive and the
 Climate Variability Diagnostics Package: New tools for the study of climate variability and
 change. *EGUsphere*, **2024**, 1–28.
- ⁹⁹⁹ McKinnon, K. A., and C. Deser, 2018: Internal variability and regional climate trends in an ¹⁰⁰⁰ observational large ensemble. *Journal of Climate*, **31** (**17**), 6783–6802.

47

- McKinnon, K. A., and C. Deser, 2021: The inherent uncertainty of precipitation variability, trends,
 and extremes due to internal variability, with implications for Western US water resources.
 Journal of Climate, 34 (24), 9605–9622.
- Merrifield, A. L., L. Brunner, R. Lorenz, V. Humphrey, and R. Knutti, 2023: Climate model
 selection by independence, performance, and spread (ClimSIPS v1. 0.1) for regional applications.
 Geoscientific Model Development, 16 (16), 4715–4747.
- Milinski, S., N. Maher, and D. Olonscheck, 2020: How large does a large ensemble need to be?
 Earth System Dynamics, **11** (**4**), 885–901.
- ¹⁰⁰⁹ Nathaniel, J., Y. Qu, T. Nguyen, S. Yu, J. Busecke, A. Grover, and P. Gentine, 2024: Chaosbench: A
- ¹⁰¹⁰ multi-channel, physics-based benchmark for subseasonal-to-seasonal climate prediction. *arXiv* ¹⁰¹¹ *preprint arXiv:2402.00712*.
- ¹⁰¹² Olonscheck, D., M. Rugenstein, and J. Marotzke, 2020: Broad consistency between observed and ¹⁰¹³ simulated trends in sea surface temperature patterns. *Geophysical Research Letters*, **47** (**10**), ¹⁰¹⁴ e2019GL086773.
- ¹⁰¹⁵ Olonscheck, D., and Coauthors, 2023: The new Max Planck Institute Grand Ensemble with CMIP6 ¹⁰¹⁶ forcing and high-frequency model output. *Journal of Advances in Modeling Earth Systems*,

¹⁰¹⁷ **15 (10)**, e2023MS003 790, https://doi.org/10.1029/2023MS003790.

- Oudar, T., P. J. Kushner, J. C. Fyfe, and M. Sigmond, 2018: No impact of anthropogenic aerosols on early 21st century global temperature trends in a large initial-condition ensemble. *Geophysical Research Letters*, **45** (17), 9245–9252.
- Penland, C., and P. D. Sardeshmukh, 1995: The optimal growth of tropical sea surface temperature
 anomalies. *Journal of climate*, 8 (8), 1999–2024.
- Po-Chedley, S., J. T. Fasullo, N. Siler, Z. M. Labe, E. A. Barnes, C. J. Bonfils, and B. D.
- ¹⁰²⁴ Santer, 2022: Internal variability and forcing influence model–satellite differences in the rate
- ¹⁰²⁵ of tropical tropospheric warming. *Proceedings of the National Academy of Sciences*, **119 (47)**, ¹⁰²⁶ e2209431119.
- Power, S., T. Casey, C. Folland, A. Colman, and V. Mehta, 1999: Inter-decadal modulation of the impact of ENSO on Australia. *Climate dynamics*, **15**, 319–324.

- Proctor, J. L., S. L. Brunton, and J. N. Kutz, 2016: Dynamic mode decomposition with control.
 SIAM Journal on Applied Dynamical Systems, **15** (1), 142–161.
- Qin, M., A. Dai, and W. Hua, 2020: Quantifying contributions of internal variability and external forcing to Atlantic multidecadal variability since 1870. *Geophysical Research Letters*, **47 (22)**, e2020GL089 504.
- Rader, J. K., C. Connolly, M. A. Fernandez, and E. M. Gordon, 2025: Attribution of the record-high
 2023 SST using a deep-learning framework. *Environmental Research Communications*.
- Rodgers, K. B., and Coauthors, 2021: Ubiquity of human-induced changes in climate variability.
 Earth System Dynamics, 12 (4), 1393–1411, https://doi.org/10.5194/esd-12-1393-2021.
- Rugenstein, M., S. Dhame, D. Olonscheck, R. J. Wills, M. Watanabe, and R. Seager, 2023:
 Connecting the SST pattern problem and the hot model problem. *Geophysical Research Letters*,
 50 (22), e2023GL105 488.
- Santer, B. D., and Coauthors, 2023: Exceptional stratospheric contribution to human finger prints on atmospheric temperature. *Proceedings of the National Academy of Sciences*, **120 (20)**,
 e2300758 120.
- Schneider, T., and I. M. Held, 2001: Discriminants of twentieth-century changes in Earth surface
 temperatures. *Journal of Climate*, 14 (3), 249–254.
- Schulzweida, U., 2023: CDO User Guide (2.3.0). Zenodo, https://doi.org/10.5281/zenodo.
 10020800.
- Seager, R., N. Henderson, and M. Cane, 2022: Persistent discrepancies between observed and
 modeled trends in the tropical Pacific Ocean. *Journal of Climate*, **35** (14), 4571–4584.
- Simpson, I. R., and Coauthors, 2025: Confronting Earth System Model trends with observations.
 Science Advances, **11** (**11**), eadt8035.
- Sippel, S., N. Meinshausen, A. Merrifield, F. Lehner, A. G. Pendergrass, E. Fischer, and R. Knutti,
 2019: Uncovering the forced climate response from a single ensemble member using statistical
 learning. *Journal of Climate*, **32** (17), 5677–5699.

- Sippel, S., N. Meinshausen, E. Székely, E. Fischer, A. G. Pendergrass, F. Lehner, and R. Knutti,
 2021: Robust detection of forced warming in the presence of potentially large climate variability.
 Science Advances, 7 (43), eabh4429.
- ¹⁰⁵⁸ Smith, D. M., and Coauthors, 2016: Role of volcanic and anthropogenic aerosols in the recent ¹⁰⁵⁹ global surface warming slowdown. *Nature Climate Change*, **6** (10), 936–940.
- Solomon, A., and M. Newman, 2012: Reconciling disparate twentieth-century Indo-Pacific ocean temperature trends in the instrumental record. *Nature Climate Change*, **2** (**9**), 691–699.
- Steinman, B. A., M. E. Mann, and S. K. Miller, 2015: Atlantic and Pacific multidecadal oscillations
 and Northern Hemisphere temperatures. *Science*, 347 (6225), 988–991.
- Stolpe, M. B., I. Medhaug, and R. Knutti, 2017: Contribution of Atlantic and Pacific multidecadal
 variability to twentieth-century temperature changes. *Journal of Climate*, **30** (16), 6279–6295.
- ¹⁰⁶⁶ Swart, N. C., and Coauthors, 2019: The Canadian Earth System Model version 5 ¹⁰⁶⁷ (CanESM5.0.3). *Geoscientific Model Development*, **12** (**11**), 4823–4873, https://doi.org/ ¹⁰⁶⁸ 10.5194/gmd-12-4823-2019.
- Tatebe, H., and Coauthors, 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.
- Taylor, K. E., 2001: Summarizing multiple aspects of model performance in a single diagram.
 Journal of geophysical research: atmospheres, **106 (D7)**, 7183–7192.
- Ting, M., Y. Kushnir, R. Seager, and C. Li, 2009: Forced and internal twentieth-century SST trends
 in the North Atlantic. *Journal of Climate*, **22** (6), 1469–1481.
- ¹⁰⁷⁶ Trenberth, K. E., and D. J. Shea, 2006: Atlantic hurricanes and natural variability in 2005. ¹⁰⁷⁷ *Geophysical research letters*, **33 (12)**.
- Varando, G., M.-Á. Fernández-Torres, J. Muñoz-Marí, and G. Camps-Valls, 2022: Learning causal
 representations with Granger PCA. UAI 2022 Workshop on Causal Representation Learning.
- Wadoux, A. M.-C., D. J. Walvoort, and D. J. Brus, 2022: An integrated approach for the evaluation
- ¹⁰⁸¹ of quantitative soil maps through Taylor and solar diagrams. *Geoderma*, **405**, 115 332.

- Watanabe, M., S. M. Kang, M. Collins, Y.-T. Hwang, S. McGregor, and M. F. Stuecker, 2024:
 Possible shift in controls of the tropical Pacific surface warming pattern. *Nature*, 630 (8016),
 315–324.
- Wills, R. C., T. Schneider, J. M. Wallace, D. S. Battisti, and D. L. Hartmann, 2018: Disentangling
 global warming, multidecadal variability, and El Niño in Pacific temperatures. *Geophysical Research Letters*, 45 (5), 2487–2496.
- Wills, R. C. J., K. C. Armour, D. S. Battisti, and D. L. Hartmann, 2019: Ocean–atmosphere dynamical coupling fundamental to the Atlantic multidecadal oscillation. *Journal of Climate*, 32 (1), 251–272.
- Wills, R. C. J., D. S. Battisti, K. C. Armour, T. Schneider, and C. Deser, 2020: Pattern recognition
 methods to separate forced responses from internal variability in climate model ensembles and
 observations. *J. Climate*, 33 (20), 8693–8719.
- Wills, R. C. J., Y. Dong, C. Proistosecu, K. C. Armour, and D. S. Battisti, 2022: Systematic climate
 model biases in the large-scale patterns of recent sea-surface temperature and sea-level pressure
 change. *Geophysical Research Letters*, 49 (17), e2022GL100011.
- Wills, R. C. J., A. L. Merrifield, A. Phillips, C. Deser, K. McKinnon, S. Po-Chedley, S. Sippel,
 and ForceSMIP Tier1 contributors, 2025: ForceSMIP Tier 1 data repository [dataset]. Zenodo,
 https://doi.org/10.5281/zenodo.15577520.
- Wyser, K., T. Koenigk, U. Fladrich, R. Fuentes-Franco, M. P. Karami, and T. Kruschke, 2021: The
 SMHI Large Ensemble (SMHI-LENS) with EC-Earth3.3.1. *Geoscientific Model Development*,
 14 (7), 4781–4796, https://doi.org/10.5194/gmd-14-4781-2021.
- ¹¹⁰³ Xu, T., M. Newman, A. Capotondi, S. Stevenson, E. Di Lorenzo, and M. A. Alexander, 2022: ¹¹⁰⁴ An increase in marine heatwaves without significant changes in surface ocean temperature ¹¹⁰⁵ variability. *Nature Communications*, **13** (1), 7396.
- ¹¹⁰⁶ Zhang, R., and T. L. Delworth, 2006: Impact of Atlantic multidecadal oscillations on India/Sahel ¹¹⁰⁷ rainfall and Atlantic hurricanes. *Geophysical research letters*, **33** (17).
- ¹¹⁰⁸ Zhang, R., and Coauthors, 2013: Have aerosols caused the observed Atlantic multidecadal variability? *Journal of the Atmospheric Sciences*, **70** (4), 1135–1144.

- ¹¹¹⁰ Zhang, X., L. Alexander, G. C. Hegerl, P. Jones, A. K. Tank, T. C. Peterson, B. Trewin, and
 ¹¹¹¹ F. W. Zwiers, 2011: Indices for monitoring changes in extremes based on daily temperature and
 ¹¹¹² precipitation data. *Wiley Interdisciplinary Reviews: Climate Change*, 2 (6), 851–870.
- Ziehn, T., and Coauthors, 2020: The Australian Earth System Model: ACCESS-ESM1.5. *JSHESS*, **70** (1), 193–214.