Forced Component Estimation Statistical Method Intercomparison Project (ForceSMIP)



<u>Robb Jnglin Wills</u>¹, Clara Deser², Karen McKinnon³, Adam Phillips², Stephen Po-Chedley⁴, Sebastian Sippel⁵, Anna Merrifield¹, and ForceSMIP contributors

¹ETH Zurich, ²NCAR, ³UCLA, ⁴LLNL, ⁵University of Leipzig

EGU 2025, CL4.10, 30th April 2025



How can we best estimate the forced and unforced components of observed climate changes, as is done with large ensembles?

Climate Trend



(Ensemble Mean) Forced Response



Internal Climate Variability





ForceSMIP aims to quantify the contributions of external forcing and internal variability in observations, using

statistical and machine learning methods

Wills et al. 2025, preprint coming soon

Climate Model

ForceSMIP Project Hackathon and Contributors

Organizers: Robert Jnglin Wills¹, Clara Deser², Karen McKinnon³, Adam Phillips², Stephen Po-Chedley⁴, Sebastian Sippel⁵

Contributors: Constantin Bône⁶, Céline Bonfils⁴, Gustau Camps-Valls⁷, Stephen Cropper³, Charlotte Connolly⁸, Shiheng Duan⁴, Homer Durand⁷, Alexander Feigin⁹, Martin Fernandez⁸, Guillaume Gastineau⁶, Andrey Gavrilov^{7,9}, Emily Gordon⁸, Moritz Günther¹⁰, Maren Höver^{11,1}, Sergey Kravtsov¹¹, Yan-Ning Kuo¹², Justin Lien¹³, Gavin Madakumbra³, Anna Merrifield¹, Nathan Mankovich⁷, Matt Newman¹⁴, Jamin Rader⁸, Jia-Rui Shi¹⁵, Sang-Ik Shin^{14,18}, Gherardo Varando⁷, Tristan Williams⁷

¹ETH Zürich, ²NCAR, ³UCLA, ⁴LLNL, ⁵University of Leipzig, ⁶LOCEAN, ⁷University of Valencia, ⁸Colorado State University, ⁹Institute of Applied Physics, RAS, ¹⁰Stanford, ¹¹Oxford, ¹²MPI-Meteorology, ¹³University of Wisconsin Milwaukee, ¹⁴Cornell, ¹⁵Tohoku University, ¹⁶NOAA PSL, ¹⁷WHOI, ¹⁸CIRES, University of Colorado















ForceSMIP Hackathon Aug. 29-31, 2023 NCAR and ETH Zurich First, some motivation: Open questions related to separating the forced response from internal variability

Are systematic discrepancies between observed and modeled SST trends forced or unforced?

Observed SST Change (1979-2020)

- No member out of 16 large ensembles simulates both the observed tropical Pacific SST gradient strengthening and the Southern Ocean cooling
- The discrepancy could be a bias in the forced response, a bias in the amplitude of multi-decadal variability, or both



Observed jet speed, SLP, and precip trends (1951-2020) are outside the model distribution. Forced or Unforced?





- CMIP6 models cannot reproduce the observed strengthening of the North Atlantic jet and related SLP and precip trends
- Knowing whether models have biases in their forced responses or internal variability is critical information for near-future water resource planning, e.g., in Southern Europe

Is observed Atlantic multi-decadal variability forced or unforced? Recent studies have come to opposite conclusions

Article

Tropical Atlantic multidecadal variability is dominated by external forcing

Chengfei He^{1⊠}, Amy C. Clement¹, Sydney M. Kramer², Mark A. Cane³, Jeremy M. Klavans², Tyler M. Fenske¹ & Lisa N. Murphy¹

Quantifying Contributions of Internal Variability and External Forcing to Atlantic Multidecadal Variability Since 1870

Minhua Qin¹ , Aiguo Dai², and Wenjian Hua¹





* multiplied by 4.82 to correct for apparent signal-to-noise problems

He et al. 2023; Qin et al. 2020

ForceSMIP: Framework, Datasets, and Statistical Methods

The ForceSMIP Challenge: Take a single realization of the climate system and estimate the forced response

1950-01



Observed (HadCRUT4) 3-month running mean surface temperature anomalies

Forced Response: Defined here as the spatiotemporally evolving anomalies (from a reference period) due to all anthropogenic and natural (volcanic, solar) external forcing

The ForceSMIP Challenge: Take a single realization of the climate system and estimate the forced response

Training data: 5 large ensembles

Evaluation data: Unlabeled individual realizations from 9 climate models (including 5 not part of the training data) and 1 observational product

• 8 field variables at monthly resolution, 1950-2022

How it works:

- 1. Participants **train** statistical methods to estimate the forced response from single realizations and apply them to the evaluation data
- 2. We **evaluate** how well these estimates reproduced the corresponding ensemble means (which were initially hidden from participants)
- 3. We examine the **observational forced response estimates** from skillful methods

How to estimate the forced response from a single realization: Methods submitted to ForceSMIP



How to estimate the forced response from a single realization: Linear methods with few tunable parameters

Low-frequency component analysis (LFCA): EOF-based method that estimates the forced response as the patterns evolving on the longest timescale (Wills et al. 2020)

Linear inverse models: Creates a statisticaldynamical model based on lead-lag information (Penland & Sardeshmukh 1995) and estimates the forced response as the pattern that evolves on the longest timescale (e.g., Frankignoul et al. 2017)

Linear regression on global mean or forcing timeseries: Estimates the forced response as anomalies covarying with a forcing or global-mean surface temperature timeseries

How to estimate the forced response from a single realization: Fingerprinting and machine learning methods

Linear fingerprinting methods:

Determine the extent to which modelbased forced response patterns show up in observations, with less weight put on noisy regions

Machine learning methods: Learns

from the training model data what a forced response pattern and/or internal variability pattern looks like and applies this information to observations (or other evaluation members), e.g., Bône et al. 2024

Evaluating the ForceSMIP methods using large ensemble evaluation data

Evaluating estimated trends (e.g., 1980-2022 SST trends)

Evaluating estimated trends (e.g., 1980-2022 SST trends)

Evaluating estimated trends with Taylor diagrams

SST Trends (1980-2022)

Precipitation Trends (1980-2022)

Note: Left plot is zoomed into the black box region of the right plot, since SST trends are more skillful

Evaluating the ForceSMIP methods: Findings

Summarizing the findings (details in the paper):

- 1. Most methods reduce RMSE in the forced response estimate compared to the raw data (for long-term and short-term trends, spatiotemporal variability, and large-scale indices)
- 2. This sometimes come at the expense of reducing the pattern correlation and/or making the amplitude to week, especially for precipitation and SLP
- 3. The relative skill of different methods varies between variables

Examining the ForceSMIP estimated forced response in observations

Wide spread of forced vs. internal attribution, even amongst methods determined to be skillful in models

ForceSMIP-mean forced response estimate differs from the forced response in climate models

Wide spread of ForceSMIP estimated forced responses in large-scale indices Methods mostly more

Conclusions and Outlook

- Many types of statistical and machine learning methods exhibit skill at estimating the forced response (removing internal variability) in single realizations
- Methods with comparable skill in the model evaluation dataset show a wide spread of forced response estimates in observations, illustrating the epistemic uncertainty in forced response estimation
 - For example, the AMV could be almost entirely forced or almost entirely unforced
- Nevertheless, ForceSMIP shows systematic differences in estimated forced responses compared to climate models
- We will release the raw ForceSMIP dataset with the forthcoming paper (*Wills et al. 2025, preprint coming soon*), and a skill-weighted observational forced response will follow after that (*Merrifield et al., in prep.*)
 - Potential applications to model evaluation, near-term climate prediction, observational large ensembles, climate variability analysis