Separation of internal and forced variability of climate using a U-Net

Constantin Bône¹, Guillaume Gastineau², Sylvie Thiria³, Patrick Gallinari⁴, and Carlos E Mejia⁵

¹LOCEAN, ISIR,Sorbonne Université/CNRS/IRD/MNHN ²LOCEAN, Sorbonne Université/CNRS/IRD/MNHN ³Sorbonne Universités, France ⁴ISIR, Sorbonne-Université ⁵LOCEAN, Sorbonne Université

August 22, 2023

Abstract

The internal variability pertains to fluctuations originating from processes inherent to the climate component and their mutual interactions. On the other hand, forced variability delineates the influence of external boundary conditions on the physical climate system. A methodology is formulated to distinguish between internal and forced variability within the surface air temperature. The noise-to-noise approach is employed for training a neural network, drawing an analogy between internal variability and image noise. A large training dataset is compiled using surface air temperature data spanning from 1901 to 2020, obtained from an ensemble of Atmosphere-Ocean General Circulation Model (AOGCM) simulations. The neural network utilized for training is a U-Net, a widely adopted convolutional network primarily designed for image segmentation. To assess performance, comparisons are made between outputs from two single-model initial-condition large ensembles (SMILEs), the ensemble mean, and the U-Net's predictions. The U-Net reduces internal variability by a factor of four, although notable discrepancies are observed at the regional scale. While demonstrating effective filtering of the El Niño Southern Oscillation, the U-Net encounters challenges in areas dominated by forced variability, such as the Arctic sea ice retreat region. This methodology holds potential for extension to other physical variables, facilitating insights into the enduring changes triggered by external forcings over the long term.

Separation of internal and forced variability of climate using a U-Net

Constantin Bône¹², Guillaume Gastineau¹, Sylvie Thiria¹, Patrick Gallinari²³and Carlos Mejia¹

5	$^1 \mathrm{UMR}$ LOCEAN, IPSL, Sorbonne Université, IRD, CNRS, MNHN
6	$^2 \mathrm{UMR}$ ISIR, Sorbonne Université, CNRS, INSERM
7	³ Criteo AI Lab

Key Points:

1

2

3

4

We present a new method to separate the forced and internal variability of the surface air temperature. We utilise a U-Net trained with global climate models outputs and implement a noise to noise methodology to eliminate internal variability. The results are assessed through the utilisation of very large ensemble simulations of two distinct climate models.

 $Corresponding \ author: \ Constantin \ B\hat{o}ne, \ {\tt constantin.bone@sorbonne-universite.fr}$

15 Abstract

The internal variability pertains to fluctuations originating from processes inherent to 16 the climate component and their mutual interactions. On the other hand, forced vari-17 ability delineates the influence of external boundary conditions on the physical climate 18 system. A methodology is formulated to distinguish between internal and forced vari-19 ability within the surface air temperature. The noise-to-noise approach is employed for 20 training a neural network, drawing an analogy between internal variability and image 21 noise. A large training dataset is compiled using surface air temperature data spanning 22 from 1901 to 2020, obtained from an ensemble of Atmosphere-Ocean General Circula-23 tion Model (AOGCM) simulations. The neural network utilized for training is a U-Net, 24 a widely adopted convolutional network primarily designed for image segmentation. To 25 assess performance, comparisons are made between outputs from two single-model initial-26 condition large ensembles (SMILEs), the ensemble mean, and the U-Net's predictions. 27 The U-Net reduces internal variability by a factor of four, although notable discrepan-28 cies are observed at the regional scale. While demonstrating effective filtering of the El 29 Niño Southern Oscillation, the U-Net encounters challenges in areas dominated by forced 30 variability, such as the Arctic sea ice retreat region. This methodology holds potential 31 for extension to other physical variables, facilitating insights into the enduring changes 32 triggered by external forcings over the long term. 33

34

Plain Language Summary

To comprehensively grasp future climate change, it becomes imperative to differ-35 entiate between forced variability and internal climate variability. Internal variability refers 36 to the climate's variations driven by the chaotic nature of geophysical fluids. Conversely, 37 forced variability denotes changes prompted by external forcings, predominantly alter-38 ations in radiative forcing, primarily due to anthropogenic activities. Here, a novel ap-39 proach is introduced for filtering internal variability through the utilisation of a convo-40 lutional neural network. This neural network is trained using a noise-to-noise method-41 ology, targeting the filtration of internal variability from surface air temperature outputs 42 of climate models or observational data. Internal variability is treated analogously to noise 43 within an image, which is removed to restore the "true image," corresponding to forced 44 variability in our case. This method capitalises on the data generated by state-of-the-45 art climate models through the coupled model intercomparison project (CMIP). To val-46

-2-

idate this methodology, we assess its performance using very large ensembles of climate
model simulations, enabling precise estimation of forced variability. Our findings demonstrate a reduction in internal variability by a factor of four, accompanied by notable regional variations.

51 **1** Introduction

The phenomenon of climate warming is characterized by an elevated surface air tem-52 perature, notably reaching a pivotal juncture during the latter half of the twentieth cen-53 tury (Eyring et al., 2021). Nevertheless, the observed anomalies in surface air temper-54 ature arise from a dual spectrum of variabilities. The first source of variability is due to 55 the effect of the external forcings, such as the increase in the greenhouse gases concen-56 tration, the variations of concentration in anthropogenic and natural aerosols, the fluc-57 tuations in solar variability or volcanic eruptions and the land-use changes. The related 58 variability is designated as the forced variability. The second source of variability is com-59 ing from processes internal to the atmosphere, oceans, cryosphere and land or the inter-60 actions between them (Cassou et al., 2018). Subsequently, this form of variability is re-61 ferred to as 'internal variability,' encapsulating its inception within the climate system 62 and its persistence even without alterations in external forcings. Despite the overarch-63 ing dominance of forced variability in shaping the broad-scale and long-term trajectory 64 of surface air temperature across the 1900-2020 timeframe (Deser et al., 2012; Kay et 65 al., 2015), a comprehensive understanding of the distinct contributions of internal and 66 forced variability remains elusive. Internal variability takes center stage in briefer tem-67 poral scales and smaller spatial dimensions. For instance, the leading mode of internal 68 variability in global air surface temperature manifests as the El Niño Southern Oscilla-69 tion (ENSO), characterized by significant anomalies in the equatorial Pacific Ocean, ac-70 companied by distant teleconnections, and a prevailing cycle spanning two to seven years 71 (Wang & Picaut, 2004). Additionally, the interdecadal Pacific variability (Newman et 72 al., 2016) and the Atlantic Multidecadal variability (Zhang et al., 2019) wield the capac-73 ity to influence climate dynamics across the decadal to multidecadal spectrum. A no-74 table example involves the deceleration in the global warming rate experienced during 75 2002-2012, commonly referred to as the global warming hiatus, which has been robustly 76 linked to Interdecadal Pacific Variability (Meehl et al., 2013; Kosaka & Xie, 2013; Eng-77 land et al., 2014). Lastly, internal variability exercises influence even over centennial and 78

-3-

multi-centennial spans (Jiang et al., 2021; S. Li & Huang, 2022) exerting substantial impact on trends within the 1900-2015 interval (Bonnet et al., 2022).

The distinction between forced variability and internal variability is essential for conducting detection and attribution studies, enabling accurate estimation and simulation of the climate's reaction to alterations in radiative forcing. Moreover, this differentiation aids in recognizing and comprehending internal climate variability. Nevertheless, the availability of instrumental observations is limited to the period since 1850, and the relatively brief duration of these observations presents challenges in effectively and confidently discerning internal variability.

For identifying both internal and forced variability, linear trends (Swart et al., 2015; 88 Vincent et al., 2015) or quadratic trends (Enfield & Cid-Serrano, 2010) have been em-89 ployed to characterize forced variability. However, linear or quadratic trends inadequately 90 capture the temporal evolution of temperature, particularly failing to account for the abrupt 91 cooling subsequent to significant volcanic eruptions, which hold significant climate im-92 pact (Schmidt et al., 2018). Additional approaches include the application of Empiri-93 cal Orthogonal Functions (EOF) analysis (Parker et al., 2007), low-frequency pattern 94 filtering (Wills et al., 2020), and linear inverse models (Marini & Frankignoul, 2014). These 95 techniques deconstruct forced variability into a combination of modes featuring distinct 96 patterns and corresponding time series. Regression analysis of the global mean surface 97 temperature (GMST) has also been employed, although this may inadvertently estab-98 lish misleading links between the Atlantic and Pacific basins (Frankignoul et al., 2017; 99 Deser & Phillips, 2023). However, a comprehensive and systematic examination of these 100 methodologies remains notably absent. 101

Climate model simulations have been employed to overcome the limitations of sparse 102 observation sampling. Conducting an ensemble of climate model simulations with diverse 103 initial conditions enables estimation of forced variability via the ensemble mean. This 104 approach effectively mitigates the variance linked to internal variability by a factor of 105 n, where n signifies the ensemble's size (Harzallah & Sadourny, 1995; Hawkins & Sut-106 ton, 2009; Ting et al., 2009; Solomon et al., 2011; Deser et al., 2014; Frankcombe et al., 107 2015). As a result, modeling centers have undertaken substantial ensembles with over 108 20 or 30 ensemble members (Jeffrey et al., 2013; Rodgers et al., 2015; Sun et al., 2018; 109 Deser et al., 2020). These large ensembles are commonly referred to as Single-Model Initial-110

-4-

Condition Large Ensembles (SMILE; Deser et al. (2020)). Multiple SMILE initiatives 111 have been undertaken using models such as CCSM3 (Collins et al., 2006), CCSM4 (Gent 112 et al., 2011), CESM (Kay et al., 2015), MPI-ESM (Maher et al., 2019), FGOALS-g3 (Li 113 et al., 2020), CanESM2 (Chylek et al., 2011), and IPSL-CM6A-LR (Bonnet et al., 2021), 114 among others. This offers a valuable dataset for crafting methodologies dedicated to the 115 disentanglement of forced and internal variability. Notably, employing members of a large 116 ensemble model as surrogate observations allows for a comparison of results with the en-117 semble mean. Differences primarily mirror residual internal variability or limitations in-118 herent in the method. 119

Nevertheless, the forced variability estimated through an ensemble mean remains 120 contingent upon the specific climate model employed. These climate models carry sub-121 stantial uncertainties, particularly in terms of their climate sensitivity (Sherwood et al., 122 2020), often attributed to factors like uncertain cloud retroaction which significantly im-123 pact equilibrium climate sensitivity (Zelinka et al., 2016). Additionally, significant un-124 certainties surround historical emissions and the linked radiative forcing from aerosols 125 (Menary et al., 2020; C. J. Smith & Forster, 2021). Moreover, the internal variability ex-126 hibited by different models also varies significantly (Parsons et al., 2020). 127

Several methodologies have been devised to harness data from diverse climate mod-128 els, as employing a multi-model approach holds the potential to alleviate the uncertain-129 ties inherent in individual climate models. Multi-model ensemble means are widely adopted 130 for estimating the forced signal (Steinman et al., 2015). Notably, techniques such as the 131 signal-to-noise-maximizing empirical orthogonal functions (Ting et al., 2009; Wills et al., 132 2020) and the discriminant analysis and maximization of the average predictability time 133 (DelSole et al., 2011) have been put forth to extract forced variability with superior ef-134 ficacy compared to ensemble means. Furthermore, scaling techniques that adjusts the 135 forced signal from models using observational data have been proposed. Among these 136 methodologies are fingerprinting methods grounded in linear regression, commonly ap-137 plied for detecting and attributing climate change with a unified forcing that encapsu-138 lates the influence of all external forcings (Hasselmann, 1993; Allen & Tett, 1999; Allen 139 & Stott, 2003). More recently, the use of scaling factors was also proposed by Frankcombe 140 et al. (2015). 141

-5-

This paper introduces an alternative approach to distinguishing internal and forced 142 variability using climate model data, employing a non-linear method that takes into ac-143 count the spatio-temporal data covariances. This method is rooted in a neural network 144 trained on data from Atmosphere-Ocean General Circulation Models (AOGCMs). Among 145 the areas where neural networks have excelled is image analysis (Egmont-Petersen et al., 146 2002). One of the prominent applications of neural networks in image processing is im-147 age denoising, involving the elimination of noise from an image to restore its true form 148 (Ilesanmi & Ilesanmi, 2021; Tian et al., 2020). In this context, internal variability is treated 149 as noise. It is demonstrated that machine learning image denoising methodologies can 150 subsequently isolate forced variability. The internal variability is eliminated, leaving be-151 hind a quantifiable residue. This method leverages the temporal and spatial information 152 inherent in climate models to establish the weights and biases of a neural network. With 153 these parameters in place, the neural network is also employed with observations to delve 154 into and attribute the progression of climate change since 1905 to 2016. To the best of 155 our knowledge, this represents the pioneering application of a dedicated neural network 156 for the purpose of disentangling internal and forced variability. 157

The structure of this paper is as follows: Section 2 outlines the data utilized. Section 3 introduces the method anchored in a neural network. Section 4 assesses the method's performance. In Section 5, the neural network method is applied to observations. Lastly, Section 6 offers the conclusion and discussion.

162 **2 Data**

163 2.1 Observations

The gridded monthly Surface Air Temperature anomaly (SAT) from 1901 to 2020, as provided by GISS Surface Temperature Analysis version 4 (GISTEMP; Hansen et al. (2010); Lenssen et al. (2019)), is employed in this study. GISTEMP amalgamates meteorological station data over land (NOAA GHCN v4) with sea surface temperature (SST) estimates from ERSST v5. This data is available on a consistent 2°x2° grid. The monthly values are aggregated to calculate annual means, and the SAT anomalies are determined using the reference period 1950-2014. 171

2.2 Climate model simulations

The monthly SAT data is sourced from historical simulations within the Coupled 172 Model Intercomparison Project Phase 5 (CMIP5; Taylor et al. (2012)) and the Coupled 173 Model Intercomparison Project Phase 6 (CMIP6; (Eyring et al., 2016)), along with sev-174 eral Single-Model Initial-Condition Large Ensembles (SMILEs) from distinct models: MPI-175 ESM (Maher et al., 2019), CSIRO-Mk3-6-0 (Collier et al., 2011), EC-Earth (Döscher et 176 al., 2021), and FGOALS-g3 (Li et al., 2020). For the historical simulations, spanning 1901 177 to 2005 (2014) for CMIP5 (CMIP6), all external forcings are integrated. These forcings 178 encompass the effects of historical greenhouse gas concentrations, anthropogenic and nat-179 ural aerosols, stratospheric ozone, solar activity, and land-use changes. Each climate model 180 delivers multiple realizations referred to as ensemble members, generated through dis-181 tinct initial conditions. From 2005 (2014 for CMIP6) until 2020, the outputs under the 182 pessimistic Representation Concentration Pathway 8.5 (RCP8.5) scenario for CMIP5 (Van Vu-183 uren et al., 2011) and the intermediate Shared Socio-economic Pathway 2 4.5 (SSP2-4.5) 184 for CMIP6 (Tebaldi et al., 2020) are employed. These simulations utilize socio-economic 185 assumptions to project future external forcing patterns. Additionally, several SMILES 186 are incorporated, employing distinct historical forcings or scenario simulations of CMIP5 187 or CMIP6 (elaborated in Table S3). While minor differences are anticipated in exter-188 nal forcing between CMIP5 and CMIP6 simulations, notable uncertainties arise in aerosol 189 emissions (C. J. Smith et al., 2020; Fyfe et al., 2021). Modest differences may also emerge 190 between the RCP8.5 (strong) and SSP2-4.5 (moderate) scenarios, particularly until 2020, 191 where actual forcings mirror observed forcings to a considerable extent (Masson-Delmotte 192 et al., 2021). 193

The count of members accessible for scenario simulations is fewer compared to the 194 historical counterparts. Therefore, we extended the outputs from historical experiments 195 using the scenario ensemble member of the same model with the same number identi-196 fication. In case the number identification is lacking, we select randomly an scenario en-197 semble member of the same climate model. 198

199

All monthly data are aggregated into annual means. Subsequently, the SAT anomalies are computed for each ensemble member using 1950-2014 as a reference period. This 200 furnishes a multi-model ensemble comprising 801 members derived from 47 AOGCMs. 201 Subsequently, the concatenated historical and scenario members are harnessed within 202

-7-

the 1901-2020 timeframe. All model data is regridded using bilinear interpolation on the horizontal grid from GISTEMP. The details pertaining to the climate model names, ensemble sizes, and the names of the employed scenario simulations are elucidated in Tabs. S1, S2, and S3.

207

2.3 Validation of the data set

The forced variability simulated within the multi-model ensemble is succinctly ex-208 amined for two specific data subsets. We investigate the MPI-ESM and FGOALS-g3 cli-209 mate models from SMILE, as they have a very large size of 100 and 115 members, re-210 spectively, which largely exceed the size of other model ensembles. Anticipatedly, the es-211 timated forced variability derived from the ensemble mean for each of these models is 212 expected to be accurate, as the reduction in variance attributed to internal variability 213 reaches 100 and 115, respectively. For instance, Deser et al. (2012, 2014) demonstrated 214 that identifying regional climate responses on time scales of several decades may neces-215 sitate between 10 to 40 members. Specifically, to detect a change in SAT between the 216 decades 2005-2014 and 2028-2037 on a global scale, the use of 3 to 6 members is requi-217 site. This requirement can surge beyond 10 for local analyses such as in North Amer-218 ica. Subsequently, the data originating from these two models is subsequently employed 219 to appraise the outcomes of the neural network model in section 4.1. 220

We utilize the ensemble mean to characterize the forced variability and employ the 221 standard deviations from the ensemble members for evaluating the internal variability. 222 Figure 1 illustrates the standard deviation of the SAT deviation from the ensemble mean 223 for FGOALS-g3 and MPI-ESM. The variability in SAT is more pronounced over land 224 surfaces ($\sim 0.3^{\circ}$ C) compared to oceans ($\sim 0.1^{\circ}$ C), consistent with the lower thermal in-225 ertia of land. Notably, substantial variability (ranging from approximately 1.5°C to 2.5°C) 226 is observed over regions coinciding with the sea ice edge, such as the Bering Sea and Nordic 227 Seas in the Northern Hemisphere, as well as the Amundsen and Weddell Seas in the South-228 ern Hemisphere. Additionally, a marked variability is observed in the equatorial Pacific 229 Ocean, with a standard deviation of 0.8°C, and this variability is more prominent in MPI-230 ESM compared to FGOALS-g3. A localized peak of variability is situated over the sub-231 polar North Atlantic, especially notable for FGOALS-g3 (reaching up to 2°C). These out-232 comes coherently reflect a significant internal variability stemming from extratropical weather 233 fluctuations over land surfaces, exhibiting local maxima around regions adjacent to the 234

-8-



Figure 1. Standard deviation of the SAT deviations from the ensemble mean for (top) MPI-ESM and (bottom) FGOALS-g3.

sea ice edge. Moreover, the variability observed in the equatorial Pacific mirrors the phenomenon of El Nino Southern Oscillation (Neelin et al., 1998).

The forced variability is estimated through the ensemble mean of each model. Sub-237 sequently, the multi-model mean (MMM) is computed by averaging the ensemble means 238 across all models, ensuring equal weight for each model. Nonetheless, MPI-ESM and FGOALS-239 g_3 are excluded from this computation, as the intention is to later compare them to the 240 MMM. To assess the prominent impact of greenhouse gas forcing, Figure 2 (a, c, e) il-241 lustrates the ensemble mean SAT anomaly for MPI-ESM, FGOALS-g3, and the MMM 242 throughout the 2010-2020 interval. Furthermore, Figure 2 (b, d, f) presents the tempo-243 ral standard deviation of the ensemble means across the period from 1901 to 2020. As 244 anticipated, all climate models project more substantial warming over land (up to 0.8°C) 245 than over oceans (approximately 0.3°C). Notably, the Arctic exhibits an amplification 246 of global warming, with warming exceeding 2°C north of 60°N. The MMM showcases an 247 average warming of 0.8°C for the 2010-2020 period, surpassing MPI-ESM (0.64°C) and 248

FGOALS-g3 (0.69°C). This aligns with the comparatively lower equilibrium climate sen-249 sitivity (ECS) of these two models (3.6°C for MPI-ESM and 2.8°C for FGOALS-g3) when 250 compared to other models employed in this study (Zelinka et al., 2020). Within the sub-251 polar Atlantic, the SAT anomalies exhibit a minimum, with negative temperatures anoma-252 lies observed in FGOALS-g3 over the Labrador Sea, or in MPI-ESM over the subpolar 253 gyre. This phenomenon, known as the North Atlantic warming hole (Keil et al., 2020), 254 is associated with a deceleration of the Atlantic meridional overturning circulation (He 255 et al., 2022). It is worth noting that such a minimum is less pronounced in the MMM, 256 presumably due to considerable uncertainties regarding the precise location of this warm-257 ing hole and the linked processes. An equivalent spatial pattern can be derived using stan-258 dard deviations, revealing values of approximately 0.3°C for the majority of global re-259 gions and higher values over land ($\sim 0.6^{\circ}$ C). Grid points located north of 60° also exhibit 260 elevated values, peaking at around 2°C in the Barents Sea for MPI-ESM or the Labrador 261 Sea for FGOALS-g3. 262

The forced variability exhibited by MPI-ESM and FGOALS-g3 diverges from that 263 of the MMM, revealing a comparatively weaker global warming trend and standard de-264 viation pattern. This divergence is particularly evident north of 60°N, where the warm-265 ing exhibits greater amplification (refer to Fig. 2), amounting to 1.54°C for MPI-ESM 266 and 1.45°C for FGOALS-g3. Local variations are also observed in regions such as the Labrador 267 Sea, Barents and Kara Sea, the Canadian archipelago, and the Bering Sea in the case 268 of FGOALS-g3. Notably, MPI-ESM similarly presents notable differences in the Barents 269 Sea. These discrepancies may arise from biases related to sea ice representation. Specif-270 ically, FGOALS-g3 depicts an excessive extent of Arctic sea ice (Li et al., 2020), which 271 in turn leads to inaccuracies in simulating the location of the sea ice edge. This discrep-272 ancy can account for spurious SAT variability attributed to the misplaced sea ice edge 273 within the Labrador Sea. The mean standard deviation of the ensemble mean registers 274 as 0.34°C for MPI-ESM and 0.43°C for FGOALS-g3, exceeding the mean standard de-275 viation of the SAT deviations of the members to the ensemble mean which is of 0.51°C 276 for MPI-ESM and 0.46°C for FGOALS-g3. This underscores that the internal variabil-277 ity is marginally more pronounced than the forced variability. 278

-10-



Figure 2. a) Ensemble mean of the air surface temperature (°C) in MPI-ESM in 2010-2020. c) Same as a) but for FGOALS-g3. e) Same as a) but for the MMM. b) Standard deviation of the ensemble mean surface air temperature (°C) in 1901-2020 for MPI-ESM. d) Same as b) but for FGOALS-g3. f) Same as b) but for the MMM.

279 3 Methods

3.1 Neural network

We design a neural network to remove the internal variability from the SAT. The 281 input data is structured with dimensions (120, 90, 180), corresponding to time spanning 282 from 1901 to 2020, latitude, and longitude, respectively. On the other hand, the output 283 holds dimensions of (112, 90, 180), encompassing the years 1905 to 2016, while maintain-284 ing the latitude and longitude dimensions intact. Notably, the output's temporal span 285 is truncated compared to the input, by excluding the initial and final four years. This 286 reduction addresses the substantial uncertainty typically observed at the dataset's end-287 points, an aspect that will be elaborated upon later. 288

A neural network's characteristics are shaped by its hyperparameters, which dictate both its architecture and training process. Our approach involves utilizing three distinct datasets, each composed of input and desired output pairs. The training dataset serves the purpose of establishing the neural network's weights and biases. Meanwhile, the validation dataset comes into play for estimating the hyperparameters. Finally, the test dataset is employed to assess the neural network's performance.

295

3.2 Constitution of the database

To construct the training dataset, we adapt a noise-to-noise methodology originally 296 introduced in Lehtinen et al. (2018). This approach was initially designed to train a neu-297 ral network in denoising images. In this method, the network is exclusively trained on 298 noisy images depicting various objects. Each object has more than one noised image de-299 picting it. In the noise to noise method, we create an input/output training database 300 that comprises pairs of noisy image combinations for identical objects. It's essential to 301 note that the network cannot effectively learn to transform a random noise realization 302 into another. Instead, the configuration is designed to approximate the mathematical 303 expectation of all noisy images associated with the same object, culminating in an es-304 timate that closely resembles the noise-free image. 305

For our application, we consider the forced spatio-temporal SAT anomalies from each climate model as distinct objects. These anomalies, inherent to each member, can be likened to noisy images, where the internal variability introduces the noise component. The ensemble members' mathematical expectation equates to the forced variabil-

ity, which can be approximated through the ensemble mean.

To create the training dataset, we follow a procedure wherein we compute pairs of 311 members for each climate model, except for MPI-ESM, FGOALS-g3, and MIROC6, which 312 are reserved for testing and validation purposes. Adopting an approach similar to Lehtinen 313 et al. (2018), we augment the dataset by introducing the ensemble mean of the climate 314 model's members as an additional member. This inclusion serves to expedite the train-315 ing process without introducing any other influences. In this process, each pair of mem-316 bers becomes an input/output pair. If we denote the number of ensemble members ob-317 tained from a specific climate model as n, this approach yields n(n+1) input/output 318 pairs per model. By accumulating such pairs from all models, the resulting training dataset 319 primarily comprises simulations characterized by the most extensive ensemble sizes (namely 320 IPSL-CM6A-LR, CanESM5, CNRM-CM6-1, and ACCESS-ESM1-5). 321

To create the validation set, we employ the ensemble simulation data from the MIROC6 model, which ranks as the third-largest ensemble in terms of size (with n = 50 members). For this purpose, we designate the ensemble members as inputs, while the ensemble mean spanning the period from 1905 to 2016 serves as the desired output.

To form the test dataset, we draw upon data derived from the FGOALS-g3 and MPI-ESM models, leveraging their extensive ensemble sizes of n = 110 and n = 100respectively. Subsequently, we proceed to make comparisons between the outputs of the neural network obtained from ensemble members and their corresponding ensemble means for both of these models.

The conclusions drawn from these tests and validation processes may exhibit some dependence on the specific model being analyzed, as alternative models could yield varying outcomes. Nevertheless, this approach has been chosen due to its simplicity and its potential to mitigate the impact of any remaining internal variability.

335 **3.3 U-Net**

Convolutional neural networks (CNNs, Yamashita et al. (2018)) constitute a category of non-linear neural networks, notably applied in tasks related to imagery (O'Shea & Nash, 2015). A distinctive attribute of CNNs is their utilization of convolutional layers, which incorporate a trainable kernel that slides across the input data.

-13-



Figure 3. Schematic of the U-Net. The arrows represent the operations within the network. The numbers shows the dimension of the data and the number of filters used.

In this context, a U-Net architecture is employed, which falls within the realm of 340 CNNs. Originally introduced by Ronneberger et al. (2015) for image segmentation, the 341 U-Net structure has gained widespread popularity in image-related analyses such as de-342 noising (Ilesanmi & Ilesanmi, 2021; Tian et al., 2020). The U-Net architecture is char-343 acterized by its inclusion of a contracting path and an expansive path, which collectively 344 give rise to its characteristic U shape (refer to Fig. 3). The contracting path adheres to 345 a conventional design of a convolutional network, featuring numerous convolutional lay-346 ers, each followed by an activation function and a max-pooling operation. As the con-347 tracting path advances, spatial information is diminished while feature information is 348 enriched. Conversely, the expansive path amalgamates feature and spatial information 349 through a sequence of up-convolutions and concatenations with high-resolution features 350 derived from the contracting path. 351

The U-Net architecture employed in this study shares similarities with the design proposed by Ronneberger et al. (2015). However, a modification is made by replacing the 2-dimensional convolutional layers with 3-dimensional counterparts. This alteration is introduced to encompass not only the spatial dimension but also the temporal dimension of the data. The selected activation function is the hyperbolic tangent. Additionally, adaptations have been made to the output layer to accommodate an output com-

-14-

prising 112 time steps. The neural network is comprised of a total of 5,659,009 trainable
 parameters.

A batch size of 8 is chosen, and the optimization process employs the Adam op-360 timizer with a learning rate of 0.001. To ensure proper application of the CNN to the 361 data, padding is introduced. This involves extending the image by appending zero val-362 ues at its edges. For the longitudinal dimension, which is periodic, the zero padding only 363 results in a slight discontinuity at 180°E, the edge of the data. Indeed, due to the na-364 ture of convolutional layers, a U-Net has more difficulty processing information located 365 at the edge of the data. This is the reason why we excluded the initial and final four years 366 (1901-1904 and 2017-2020) in the U-Net's outputs. The chosen cost function is the root 367 mean squared error (RSME), calculated using an area-weighted mean of the gridded data. 368

The validation dataset is utilized to determine the optimal values for two key hy-369 perparameters: the number of epochs and the number of filters used in the convolutional 370 layers. The term "number of filters" pertains to the thickness of the convolutional lay-371 ers. The number of epochs refers to how many times the training dataset is processed 372 during the training phase. These hyperparameters are selected to minimize the root mean 373 squared error (RMSE) using the validation dataset. Examination of the validation RMSE 374 for different values of epochs and layer thickness reveals a consistent pattern (see Fig. 375 S1): a significant reduction in RMSE occurs in the initial epochs, followed by a grad-376 ual increase. As a result, we settle on a layer thickness of 16 for the first layer (as shown 377 in Fig. 3) and a total of 32 epochs. 378

379 **3.4 Example**

Figure 4 provides an illustrative example featuring two randomly selected ensem-380 ble members from MPI-ESM and FGOALS-g3. The comparison focuses on the SAT at 381 the year 2016, depicted in the top panels, as well as the resulting output generated by 382 the neural network in 2016 (centre panels), juxtaposed against the ensemble mean anomaly 383 for the same year (bottom panels). The anticipated impact of elevated greenhouse gas 384 concentrations in 2016 is evident in the SAT of both MPI-ESM and FGOALS-g3 mem-385 bers, which exhibit warm anomalies. However, the internal variability introduces anoma-386 lies that surpass those of the ensemble mean in numerous regions, accompanied by some 387 negative anomalies in other areas. To elaborate, an instance of cooling is simulated across 388

-15-

the Equatorial Pacific Ocean, possibly linked to a La Niña event in the case of MPI-ESM. 389 The same ensemble member displays cooling over land in equatorial Africa, South-Eastern 390 Asia, and Australia, as well as in extratropical zones like the North Atlantic Ocean and 391 the Weddell Sea. In the example from FGOALS-g3, cold anomalies emerge over the Nordic 392 Seas and the Labrador Sea. Such cooling diverges from the ensemble average, which ex-393 hibits a relatively uniform warming pattern across the globe, with a more pronounced 394 effect over landmasses. Notably, the Arctic and its environs experience heightened warm-395 ing compared to other global regions, due to polar amplification. Conversely, minimal 396 warming is observed in the Southern Ocean and the subpolar North Atlantic Ocean, and 397 even a cooling tendency is noted in the Northern Atlantic warming hole. 398

The SAT obtained from the U-Net's output, utilizing the same ensemble member 399 as input, exhibits a pattern strikingly similar to that of the ensemble mean (compare cen-400 tre and bottom panels). In both instances, the pattern is relatively uniform, albeit with 401 heightened warming observed over land areas, coupled with an Arctic Amplification phe-402 nomenon. This suggests that the internal variability—such as the influence of ENSO events 403 or the effects of prolonged weather patterns over continents—has been successfully elim-404 inated. The regions displaying subdued warming or cooling tendencies are replicated, 405 although the exact positioning and intensity might not precisely match those of the en-406 semble mean in certain areas, particularly the Southern and subpolar North Atlantic. 407 It's worth noting a minor discontinuity at 180°E resulting from the padding process. 408

The performance of the method is quantified more systematically in the next sec-tion.

411

4 The U-Net as an internal variability filter

The U-Net was applied to every member of FGOALS-g3 and MPI-ESM. We then compare the results obtained with the respective ensemble mean of these two climate models.

Figures 5a and 5b illustrate the root mean squared error (RMSE) between the outcomes generated by the U-Net and the corresponding ensemble mean for the time period of 1905-2016. Notably, the discrepancies in U-Net's predictions are not uniformly distributed across space. The RMSE values fall within the range of 0.05°C to 0.5°C. The discrepancies generally remain below 0.2°C in tropical regions, except for instances over

-16-



Figure 4. (First column) Anomalies of SAT in a randomly chosen member of MPI-ESM, the associated U-Net output and ensemble mean in 2016. (Second column) Same as the first column but for a randomly chosen ensemble member for FGOALS-g3.

Western Africa in the MPI-ESM model. In contrast, the largest errors are concentrated 420 in polar areas, encompassing the Nordic Seas, Labrador Sea, and Bering Sea. Moreover, 421 sizable errors are also evident over the Southern Ocean and the continents of the North-422 ern Hemisphere situated above 45°N. These high-error regions correspond to locales char-423 acterized by substantial internal variability (refer to Figure 1). Nevertheless, it is note-424 worthy that the errors produced by the U-Net are approximately five times smaller than 425 the actual internal variability. Between the years 1996 and 2016, both ensemble results 426 exhibit a warming trend that is roughly 0.1°C lower in the U-Net results when compared 427 to the ensemble mean (as observed in Figs. 5cd). This difference is indicated by the nearly 428 consistent negative divergence situated between latitudes 45°N and 45°S. 429

The prevailing trend of systematic underestimation is, however, disrupted by an 430 exception involving the subpolar Atlantic and the Southern Ocean, where an overesti-431 mation of warming is observed. This overestimation is particularly conspicuous in the 432 FGOALS-g3 model, with warming anomalies extending to approximately 1°C over the 433 Labrador Sea and 0.5°C over the Bering Sea. This divergence from the ensemble mean 434 highlights the limited capacity of the neural network to accurately predict forced changes 435 within the subpolar North Atlantic, which is a region that exhibits inconsistent surface 436 temperature shifts across models (Swingedouw et al., 2021). The neural network's per-437 formance is restricted due to this discrepancy among models, which hampers its abil-438 ity to discern the specific features of each climate model. For example, in the case of FGOALS-439 g3, the extensive anomalies in the Labrador and Bering Seas are not mirrored in the multi-440 model mean (see Figure 2). It's also plausible that the substantial internal variability 441 observed in these regions poses a challenge for accurate removal by the neural network 442 (refer to Figure 1). This underestimation extends to the continents, with a greater im-443 pact on South America, Africa, and Australia in the tropics, as well as North America 444 and Northern Siberia in boreal regions. The degree of underestimation reaches 0.15°C 445 for MPI-ESM and 0.13°C for FGOALS-g3 in these regions. 446

Figures 6c and 6d illustrate the temporal evolution of the global surface air temperature (GSAT) for both the MPI-ESM and FGOALS-g3 models, before and after applying the U-Net correction. The range of data variability is portrayed by a 90% confidence interval assuming an Gaussian distribution. The forced variability's temporal trend extracted via ensemble mean (depicted by the red line) is effectively captured by the U-Net outputs (represented by the blue line and blue shading).

-18-



Figure 5. a) Root mean square difference of the surface air temperature, in C°, between the outputs of the U-Net and the mean ensemble in MPI-ESM, calculated across the members and all years in 1905-2016 b. b) Same as a) but for FGOALS-g3 c) Difference of the time mean SAT anomaly during 1996-2016, in °C, between the mean output of the U-Net and the corresponding ensemble mean, for MPI-ESM. d) Same as c) but for FGOALS-g3

- From 1905 to 2016, a GSAT rise is observed, aligning with the anticipated shifts in radiative forcing (Gulev et al., 2021). Additionally, a cooling pattern emerges a few years subsequent to the significant volcanic eruptions of Agung (1963), El Chichón (1982), and Pinatubo (1991), a phenomenon accurately estimated by the U-Net. This outcome aligns with expectations based on climate models incorporating volcanic aerosol emissions. Impressively, the U-Net's outputs exhibit a marginal spread, reduced approximately tenfold, indicating a substantial removal of internal variability.
- Nonetheless, the U-Net results exhibit anomalies with a slightly diminished am-460 plitude compared to the ensemble mean. The spread of the U-Net outputs is also ap-461 proximately twice as wide at the time series' beginning and end. The distribution of spa-462 tially averaged RMSE values within 90°S-90°N, comparing all U-Net outputs to the en-463 semble mean (depicted in Fig. 6a and 6b as blue histograms), reveals errors of around 464 0.12°C in MPI-ESM and 0.13°C in FGOALS-g3. Additionally, we examine the RMSE 465 values when averaging within 60°N-90°N, as Fig. 5ab suggests that errors are most pro-466 nounced in this region (illustrated in Fig. 6ab as red histograms). Errors north of 60°N 467

are approximately twice as substantial as global averages, with an average error of around 468 0.23°C in MPI-ESM and 0.26°C in FGOALS-g3. In Fig. 6ef, the internal variability ob-469 served when averaging the SAT north of $60^{\circ}N$ (as depicted by the red shading) is con-470 siderable in the raw model outputs (around 0.8° C). The ensemble mean SAT anomalies 471 in this region increase from approximately -1°C in the early twentieth century to about 472 1.2°C in 2010. The temporal evolution of the SAT north of 60°N demonstrates notable 473 similarity between the ensemble mean and the ensemble mean of U-Net outputs, with 474 a roughly 10-fold reduction in spread. However, the amplitude of the anomalies is slightly 475 underestimated, with a reduction of around 0.3°C in negative anomalies in the U-Net 476 output between 1905 and 1930 in MPI-ESM. For FGOALS-g3, the SAT is underestimated 477 by around $0.2^\circ\mathrm{C}$ during 1970-1990. 478

In Figure S2, the quadratic errors between the mean ensemble members and the 479 U-Net output are presented for each year, with global (90°S-90°N) and north of 60°N av-480 erages considered for both MPI-ESM and FGOALS-g3. Notably, the RMSE exhibits el-481 evated values during the initial and final years, characterized by peaks around the years 482 1975-1985 in both models. This pattern underscores the presence of substantial uncer-483 tainties at the data's onset and conclusion. When applying the 1900-2020 period for the 484 output (without excluding the first and last four years), the errors actually surpass those 485 portrayed in Figure S2, a fact that elucidates the rationale for excluding the endpoints 486 in the ongoing analysis, as detailed in the methods (section 2). Moreover, the notable 487 error peak during 1975-1985 lacks a definitive explanation, although it's plausible that 488 this discrepancy could be linked to uncertainties associated with the implementation of 489 aerosol forcings, notably CMIP5 for MPI-ESM and CMIP6 for FGOALS-g3. 490

The errors exhibited by the U-Net in relation to data from FGOALS-g3 are more prominent compared to those arising from the use of MPI-ESM data. This discrepancy can be attributed to the fact that MPI-ESM's simulated forced variability aligns more closely with the training data's characteristics, on average. Specifically, the training data's forced variability is in line with that of the MMM, and MPI-ESM demonstrates a smaller root mean squared difference from the MMM compared to FGOALS-g3 (as illustrated in Fig. 2).

To assess the reduction in internal variability achieved by the U-Net, we can quantitatively measure the number of ensemble members needed to surpass the U-Net's in-

-20-



Figure 6. a) Histogram showing the distribution of the RMSE between the mean ensemble and the U-Net outputs of MPI-ESM. b) Same as a), but for FGOALS-g3. c) Time evolutions of the global mean surface air temperature, in °C, for the ensemble mean and the mean U-Net outputs for MPI-ESM. Color shade shows the spread of the time series, with 90% the ensemble members uncertainty assuming a gaussian distribution. d) Same as c) but for FGOALS-g3. e) and f) are the same as c) and d) but when averaging the SAT, in °C, north of 60°N.

dividual member results using a basic ensemble mean approach. This evaluation is con-500 ducted through a random subsampling process involving 500 sets of m members, where 501 m varies from 1 to 40, for both the FGOALS-g3 and MPI-ESM ensembles. Within each 502 subset, ensemble means are calculated. The RMSE between these subsample ensemble 503 means and the actual ensemble mean obtained from all members is then determined (de-504 picted by vertical red and blue lines in Figure 7). This RMSE computation is performed 505 across all grid points and is spatially averaged. The 90% intervals, assuming an Gaus-506 sian distribution, of the 500 subsamples are also illustrated. This analysis is done for both 507 the MPI-ESM and FGOALS-g3 ensembles across distinct geographical regions: global 508 (90°S-90°N), North Atlantic (60°W-0°E, 0°N-60°N), North Pacific (120°E-100°W, 20°N-509 60° N), Niño3 (5°N-5°S, 150°W-90°W), as well as polar regions north of 60°N and south 510 of 60°S. These chosen regions exhibit considerable forced and internal variability, as vi-511 sually demonstrated in Fig. 1 and Fig. 2. Additionally, this evaluation is extended to 512 encompass both oceanic and terrestrial areas in the 60°S-60°N band, allowing for a more 513 comprehensive understanding of the U-Net's performance. The horizontal lines in the 514 illustration correspond to the same RMSE values but for the U-Net output from each 515 individual member. The accompanying color shade represents the spread of 90% uncer-516 tainty assuming an Gaussian distribution. 517

Figure 7a visually illustrates the progression of errors within the subset of mem-518 bers as the size of the subset increases. This pattern aligns with expectations, as a larger 519 subset size leads to better estimations of forced variability and a corresponding reduc-520 tion in residual internal variability by a factor of \sqrt{n} . The distribution of U-Net outputs 521 mirrors the histograms presented in Figure 6, showing a high degree of similarity across 522 both climate models. The U-Net effectively diminishes internal variability in GSAT by 523 approximately a factor of slightly more than four, which is analogous to the residual vari-524 ability observed within subsets containing around 17 members for FGOALS-g3 and 20 525 members for MPI-ESM. When focusing on regions spanning oceans and land between 526 60°N and 60°S, the outcomes remain largely consistent, showcasing a reduction in error 527 magnitude by a factor of approximately four. This reduction corresponds closely to that 528 achieved by using a subset of 15 to 20 members. 529

The U-Net's efficacy stands out prominently over the equatorial Pacific region, as depicted in panel 7f. This region is known for being heavily influenced by the ENSO, which dominates internal variability. The U-Net achieves a substantial reduction in variabil-

-22-

ity, amounting to a factor of 5.5. This reduction is akin to the outcome of utilizing an
 ensemble mean derived from around 30 members for both MPI-ESM and FGOALS-G3.

In other regions, the variability reduction is quite similar to that found globally. 535 For instance, this consistency is observed in the North Pacific and polar regions, where 536 the required number of members for equivalent outcomes remains relatively steady. How-537 ever, in terms of removing internal variability, the U-Net showcases higher efficiency in 538 the context of MPI-ESM for most scenarios. This pattern holds true except for the North 539 Atlantic, where a notable deviation is observed: a set of 15 members is necessary in MPI-540 ESM to achieve results equivalent to the U-Net (~4-fold reduction in residual variabil-541 ity), while merely 5 members suffice for FGOALS-g3 (halving of the residual variabil-542 ity). 543

The variation in performance between FGOALS-g3 and MPI-ESM might arise from dissimilarities in their internal variability, particularly over multi-decadal timescales, or due to differences in forced variability compared to the training data. Having completed this method evaluation, our focus now shifts to examining the outcomes when the U-Net is employed with observational data.

549

4.1 Filtering of the observations

The U-Net is now employed to process SAT observations derived from GISSTEMP. 550 By utilizing observed data as input, the U-Net provides an estimate of the forced vari-551 ability. In the interval from 1996 to 2016, the U-Net-derived forced SAT (depicted in Fig-552 ure 8a) illustrates a fairly uniform warming, with amplified warming evident over the 553 Arctic region, consistent with Arctic amplification. Furthermore, this warming effect is 554 slightly more pronounced over land compared to oceans. Conversely, the Southern Ocean 555 experiences less warming in comparison to other global regions. The spatial distribution 556 of standard deviations (Figure 8b), computed from 1905 to 2016 using U-Net output, 557 mirrors the anomalies observed in the 1996-2016 period. This agreement indicates the 558 prevailing influence of increasing anthropogenic forcing. Notably, this pattern closely re-559 sembles the changes observed in the multi-model mean (MMM) (as depicted in Fig. 2). 560 This underscore the significant contribution of the training dataset in determining the 561 identified forced changes. 562

-23-



Figure 7. Spatial average of the RMSE for the forced variability estimated with the U-Net outputs obtained from each ensemble member, and the forced variability obtained with ensemble averages subsampling ensemble of size 1 to 40; for (red) MPI-ESM and (blue) FGOALS-g3. The RMSE calculated from the U-Net and each ensemble member is given by (color shade) the interval including 90% of the distribution, assuming a gaussian distribution, and (horizontal dashed line) the mean RMSE. The RMSE calculated from 500 subsample of size between 1 to 40 is illustrated with (vertical lines) the intervals including 90% of the ensemble member distribution, also assuming a gaussian distribution.



Figure 8. Forced surface air temperature (in °C) anomaly when applying the U-Net to GIS-STEMP observation : a) time average in 1996-2016; b) standard deviation in 1905-2016.



Figure 9. Standard deviation of the SAT deviations from the forced SAT, as estimated using the U-Net, in 1905-2016.

To quantify internal variability within the observations, we compute the deviations 563 of observed SAT anomalies from the estimated forced changes. The resulting internal 564 variability pattern, illustrated by the time standard deviation of these deviations shown 565 in Figure 9, mirrors the model-derived pattern (Fig. 1). Higher internal variability val-566 ues are observed over land areas, as well as regions near the boundaries of sea ice, such 567 as the Labrador Sea and the Nordic Seas in the Northern Hemisphere, and the South-568 ern Ocean. Notably, a local maximum of internal variability emerges in the equatorial 569 Pacific, corresponding to the El Niño-Southern Oscillation region. This similarity in the 570 spatial distribution of internal variability between observations and models underscores 571 the consistency of our findings. 572

We now shift our focus to the GSAT and the Niño 3.4 region (5°N-5°S, 170°W-120°W), 573 with a particular emphasis on Niño 3.4 due to its notably improved performance in our 574 study. In the global context (Figure 10a), the forced variability reveals a consistent warm-575 ing trend, which becomes more pronounced during the 1960s. Notably, the major vol-576 canic eruptions of Agung (1963), El Chichón (1982), and Pinatubo (1991) are associated 577 with temporary cooling patterns. By 2016, the GSAT anomaly reaches 0.7°C. As expected, 578 the forced variability time series exhibits a significant reduction in inter-annual variabil-579 ity. This reduction is particularly striking within the Niño 3.4 region (Figure 10b), where 580 variability at 2 to 7 years is almost entirely eliminated. The U-Net estimates the Niño 581

-26-



Figure 10. Time series of (red) the observed SAT anomaly and (blue) the forced SAT anomaly estimated by the U-Net for a) the global mean b) NINO 3.4 and c) the relative SAT, calculated as the difference between the averaged SAT in Niño 3.4 region and the tropical ocean SAT (30°S-30°N).

3.4 forced variability, depicting a steady warming trend. To quantify the changes of SAT 582 in Niño 3.4 relative to the tropics, we calculate also the relative SAT, defined as the dif-583 ference between the average SAT on the NINO 3.4 region and the average SAT on ocean 584 grid between 30°S-30°N. The relative SST shows that the warming over the Niño 3.4 fol-585 lows that of the tropics, so that no clear El Niño-like reponse is found, unlike climate 586 models (Fig. 2). Some authors (Clement et al., 1996; Heede et al., 2020) have suggested 587 that a forced cooling could exists in the relative SAT, called thermostat effect. Here the 588 relative SAT shows a very small cooling (see Fig. 10c). In addition the SAT in the Niño 589 3.4 region are not affected by the forcing from the main volcanic eruptions. Therefore, 590 no evidence of a Niño-like response to volcanic eruption (as in Khodri et al. (2017)) is 591 found.

592

593 5 Conclusion

A novel approach is introduced in this study to effectively eliminate internal vari-594 ability from a time-evolving two-dimensional dataset, specifically focusing on surface air 595 temperature. The method employs a U-Net neural network and draws inspiration from 596 the noise-to-noise technique. This framework treats internal variability as an analogous 597 noise superimposed on the underlying forced variability. The U-Net model is trained us-598 ing outputs from a diverse ensemble of climate models obtained from the CMIP simu-599 lations. Subsequently, this trained network is applied to observational data to unveil the 600 forced variability signal by attenuating internal variability. The validation of this method 601 involves utilizing large ensemble simulations from individual models, specifically the MPI-602 ESM and FGOALS-g3, to gauge its effectiveness. The forced variability derived from the 603 ensemble mean is then contrasted with the outcomes from the U-Net application. To quan-604 titatively assess the U-Net's efficacy in reducing internal variability, an "equivalent en-605 semble size" is computed. This metric indicates the ensemble size that would be required 606 to achieve the same level of precision in capturing forced changes as the U-Net which is 607 applied to a single member. The U-Net outputs for these two climate models' test data 608 exhibit an error equivalent to an internal variability reduction of a factor of more than 609 4. This magnitude corresponds to the internal variability one could expect from an en-610 semble averaging 17 to 20 members. Furthermore, when the U-Net is applied to surface 611 air temperature observations, the inferred forced changes align closely with the multi-612 model mean in terms of spatial patterns. The U-Net's results do not suggest an El Niño-613 like response to global warming. We observe that the U-Net encounters greater challenges 614 in accurately estimating forced variability over the Arctic region. This discrepancy can 615 be attributed to the significant forced and internal variability associated with changes 616 in sea-ice extent in that area. Additionally, the U-Net's performance in capturing forced 617 variability in the North Atlantic is less successful for the FGOALS-g3 model. This lim-618 itation might be linked to uncertainties stemming from the multi-decadal variability preva-619 lent in these regions (Menary & Wood, 2018; Zhang, 2007). 620

621 622

623

624

625

In the pursuit of enhancing the U-Net methodology, several avenues for future improvements have been identified. One potential approach is to address the U-Net's sensitivity to the multi-model consensus of future variability by employing neural network regularization techniques, such as weights penalisation. Additionally, preprocessing methods like data augmentation could be explored to potentially mitigate such impacts. Im-

-28-

proving the evaluation process of the U-Net's performance is also on the horizon. This 626 could involve testing the U-Net on a broader range of climate models to assess its gen-627 eralizability. Comparing its outcomes with results from alternative methods, such as signal-628 to-noise filtering, could offer a comprehensive evaluation of the U-Net's effectiveness. To 629 broaden the scope of application, the U-Net's performance might be further investigated 630 using additional climate variables beyond surface air temperature (SAT). Variables such 631 as sea level surface pressure and precipitation could be explored, capitalizing on poten-632 tial correlations among these variables to provide more comprehensive insights. Lastly, 633 the proposed method holds the potential for wider applications, including its deployment 634 on simulations from projects like the Detection and Attribution Model Intercomparison 635 Project (Gillett et al., 2016) or the Large Ensemble Single Forcing Model Intercompara-636 ison Project (D. M. Smith et al., 2022). By leveraging transfer learning, the U-Net trained 637 on historical simulations could be adapted to these datasets. This adaptation could fa-638 cilitate the evaluation of specific forcing effects in individual climate models, offering a 639 valuable tool for studying the impact of different external factors on the climate system. 640 Such extensions of the method could contribute significantly to our understanding of cli-641 mate attribution and variability. 642

643 Acknowledgments

We acknowledge the support of the SCAI doctoral program managed by the ANR with
the reference ANR-20-THIA-0003, the support of the EUR IPSL Climate Graduate School
project managed by the ANR under the "Investissements d'avenir" programme with the
reference ANR-11-IDEX-0004-17-EURE-0006. This work was performed using HPC resources from GENCI-TGCC A0090107403 and A0110107403, and GENCI-IDRIS AD011013295.
Guillaume Gastineau was funded by the JPI climate/JPI ocean ROADMAP project (grant
number ANR-19-JPOC-003).

651 6 Open Research

652

Data Availability Statement

The CMIP5 and CMIP6 data is available through the Earth System Grid Feder-

ation and can be accessed through different international nodes. For example : https://

esgf-node.ipsl.upmc.fr/projects/esgf-ipsl/

656 657

677

Codes used in this article for the backward optimization and the figures are from Bône (2023) software available freely at https://zenodo.org/record/8233743.

658 References

- Allen, M. R., & Stott, P. A. (2003). Estimating signal amplitudes in optimal finger printing, Part I: Theory. *Climate Dynamics*, 21, 477–491.
- Allen, M. R., & Tett, S. F. (1999). Checking for model consistency in optimal finger printing. *Climate Dynamics*, 15, 419–434.
- Bonnet, R., Boucher, O., Deshayes, J., Gastineau, G., Hourdin, F., Mignot, J., ...
- Swingedouw, D. (2021). Presentation and evaluation of the ipsl-cm6a-lr ensemble of extended historical simulations. Journal of Advances in Modeling Earth
 Systems, 13(9), e2021MS002565.
- Bonnet, R., Boucher, O., Vrac, M., & Jin, X. (2022). Sensitivity of bias adjustment
 methods to low-frequency internal climate variability over the reference period:
 an ideal model study. *Environmental Research: Climate*, 1(1), 011001.
- Bône, C. (2023). Codes for "Separation of internal and forced variability of climate
 using a U-Net" [Software]. Retrieved from https://zenodo.org/record/
 8233743
- Cassou, C., Kushnir, Y., Hawkins, E., Pirani, A., Kucharski, F., Kang, I.-S., &
 Caltabiano, N. (2018). Decadal Climate Variability and Predictability: Challenges and Opportunities. Bulletin of the American Meteorological Society,
 99(3), 479 490. Retrieved from https://journals.ametsoc.org/view/

journals/bams/99/3/bams-d-16-0286.1.xml

- ⁶⁷⁸ Chylek, P., Li, J., Dubey, M., Wang, M., & Lesins, G. (2011). Observed and model
 ⁶⁷⁹ simulated 20th century Arctic temperature variability: Canadian earth system
 ⁶⁸⁰ model CanESM2. Atmospheric Chemistry and Physics Discussions, 11(8),
 ⁶⁸¹ 22893–22907.
- ⁶⁶² Clement, A. C., Seager, R., Cane, M. A., & Zebiak, S. E. (1996). An ocean dynami ⁶⁶³ cal thermostat. *Journal of Climate*, 9(9), 2190–2196.
- Collier, M. A., Jeffrey, S. J., Rotstayn, L. D., Wong, K., Dravitzki, S., Moseneder,
- C., ... others (2011). The CSIRO-Mk3. 6.0 Atmosphere-Ocean GCM: partici pation in CMIP5 and data publication. In *International congress on modelling and simulation-modsim* (pp. 2691–2697).

688	Collins, W. D., Rasch, P. J., Boville, B. A., Hack, J. J., McCaa, J. R., Williamson,
689	D. L., Zhang, M. (2006). The formulation and atmospheric simulation of
690	the Community Atmosphere Model version 3 (CAM3). Journal of Climate,
691	19(11), 2144-2161.
692	DelSole, T., Tippett, M. K., & Shukla, J. (2011). A significant component of un-
693	forced multidecadal variability in the recent acceleration of global warming.
694	Journal of Climate, 24(3), 909–926.
695	Deser, C., Lehner, F., Rodgers, K. B., Ault, T., Delworth, T. L., DiNezio, P. N.,
696	\dots others (2020). Insights from Earth system model initial-condition large
697	ensembles and future prospects. Nature Climate Change, $10(4)$, 277–286.
698	Deser, C., Phillips, A., Bourdette, V., & Teng, H. (2012). Uncertainty in climate
699	change projections: the role of internal variability. Climate dynamics, $38, 527-$
700	546.
701	Deser, C., & Phillips, A. S. (2023). A range of outcomes: the combined effects of
702	internal variability and anthropogenic forcing on regional climate trends over
703	Europe. Nonlinear Processes in Geophysics, $30(1)$, 63–84.
704	Deser, C., Phillips, A. S., Alexander, M. A., & Smoliak, B. V. (2014). Projecting
705	North American climate over the next 50 years: Uncertainty due to internal
706	variability. Journal of Climate, 27(6), 2271–2296.
707	Döscher, R., Acosta, M., Alessandri, A., Anthoni, P., Arneth, A., Arsouze, T.,
708	others (2021) . The EC-earth3 Earth system model for the climate model in-
709	tercomparison project 6. Geoscientific Model Development Discussions, 2021,
710	1 - 90.
711	Egmont-Petersen, M., de Ridder, D., & Handels, H. (2002). Image processing with
712	neural networks—a review. Pattern recognition, $35(10)$, 2279–2301.
713	Enfield, D. B., & Cid-Serrano, L. (2010). Secular and multidecadal warmings in the
714	North Atlantic and their relationships with major hurricane activity. Interna-
715	tional Journal of Climatology: A Journal of the Royal Meteorological Society,
716	30(2), 174-184.
717	England, M. H., McGregor, S., Spence, P., Meehl, G. A., Timmermann, A., Cai,
718	W., Santoso, A. (2014). Recent intensification of wind-driven circulation
719	in the Pacific and the ongoing warming hiatus. Nature climate change, $4(3)$,
720	222-227.

721	Eyring, V., Bony, S., Meehl, G. A., Senior, C. A., Stevens, B., Stouffer, R. J., &
722	Taylor, K. E. (2016). Overview of the Coupled Model Intercomparison Project
723	Phase 6 (CMIP6) experimental design and organization. Geoscientific Model
724	$Development, \ 9(5), \ 1937-1958.$
725	Eyring, V., Gillett, N., Rao, K. A., Barimalala, R., Parrillo, M. B., Bellouin, N.,
726	Zho, B. (2021). Human Influence on the Climate System. In Climate Change
727	2021: The Physical Science Basis. Contribution of Working Group I to the
728	Sixth Assessment Report of the Intergovernmental Panel on Climate Change.
729	Cambridge University Pres.
730	Frankcombe, L. M., England, M. H., Mann, M. E., & Steinman, B. A. (2015). Sep-
731	arating internal variability from the externally forced climate response. $Journal$
732	of Climate, $28(20)$, $8184-8202$.
733	Frankignoul, C., Gastineau, G., & Kwon, YO. (2017). Estimation of the SST
734	response to anthropogenic and external forcing and its impact on the Atlantic
735	multidecadal oscillation and the Pacific decadal oscillation. Journal of Climate,
736	30(24), 9871-9895.
737	Fyfe, J. C., Kharin, V. V., Santer, B. D., Cole, J. N., & Gillett, N. P. (2021). Sig-
738	nificant impact of forcing uncertainty in a large ensemble of climate model
739	simulations. Proceedings of the National Academy of Sciences, 118(23),
740	e2016549118.
741	Gent, P. R., Danabasoglu, G., Donner, L. J., Holland, M. M., Hunke, E. C., Jayne,
742	S. R., others (2011). The community climate system model version 4.
743	Journal of climate, 24(19), 4973–4991.
744	Gillett, N. P., Shiogama, H., Funke, B., Hegerl, G., Knutti, R., Matthes, K.,
745	Tebaldi, C. (2016). The detection and attribution model intercomparison
746	project (DAMIP v1. 0) contribution to CMIP6. Geoscientific Model Develop-
747	$ment, \ 9(10), \ 3685-3697.$
748	Gulev, S. K., Thorne, P. W., Ahn, J., Dentener, F. J., Domingues, C. M., Gerland,
749	S., others (2021). Changing state of the climate system.
750	Hansen, J., Ruedy, R., Sato, M., & Lo, K. (2010). Global surface temperature
751	change. Reviews of Geophysics, $48(4)$.
752	Harzallah, A., & Sadourny, R. (1995). Internal versus SST-forced atmospheric vari-
753	ability as simulated by an atmospheric general circulation model. Journal of

754	$Climate, \ 8(3), \ 474-495.$
755	Hasselmann, K. (1993). Optimal fingerprints for the detection of time-dependent cli-
756	mate change. Journal of Climate, $6(10)$, 1957–1971.
757	Hawkins, E., & Sutton, R. (2009). The potential to narrow uncertainty in regional
758	climate predictions. Bulletin of the American Meteorological Society, $90(8)$,
759	1095-1108.
760	He, C., Clement, A. C., Cane, M. A., Murphy, L. N., Klavans, J. M., & Fenske,
761	T. M. (2022). A North Atlantic warming hole without ocean circulation.
762	Geophysical research letters, $49(19)$, e2022GL100420.
763	Heede, U. K., Fedorov, A. V., & Burls, N. J. (2020). Time scales and mechanisms
764	for the tropical Pacific response to global warming: A tug of war between the
765	ocean thermostat and weaker Walker. Journal of Climate, $33(14)$, $6101-6118$.
766	Ilesanmi, A. E., & Ilesanmi, T. O. (2021). Methods for image denoising using convo-
767	lutional neural network: a review. Complex & Intelligent Systems, $7(5)$, 2179–
768	2198.
769	Jeffrey, S., Rotstayn, L., Collier, M., Dravitzki, S., Hamalainen, C., Moeseneder, C.,
770	\ldots Syktus, J. (2013). Australia's CMIP5 submission using the CSIRO-Mk3. 6
771	model. Australian Meteorological and Oceanographic Journal, $63(1)$, 1–13.
772	Jiang, W., Gastineau, G., & Codron, F. (2021). Multicentennial variability driven
773	by salinity exchanges between the Atlantic and the Arctic Ocean in a cou-
774	pled climate model. Journal of Advances in Modeling Earth Systems, 13(3),
775	e2020MS002366.
776	Kay, J. E., Deser, C., Phillips, A., Mai, A., Hannay, C., Strand, G., others
777	(2015). The Community Earth System Model (CESM) large ensemble project:
778	A community resource for studying climate change in the presence of internal
779	climate variability. Bulletin of the American Meteorological Society, $96(8)$,
780	1333–1349.
781	Keil, P., Mauritsen, T., Jungclaus, J., Hedemann, C., Olonscheck, D., & Ghosh, R.
782	(2020). Multiple drivers of the North Atlantic warming hole. <i>Nature Climate</i>
783	$Change, \ 10(7), \ 667-671.$
784	Khodri, M., Izumo, T., Vialard, J., Janicot, S., Cassou, C., Lengaigne, M., oth-
785	ers (2017). Tropical explosive volcanic eruptions can trigger El Niño by cooling
786	tropical Africa. Nature communications, $\mathcal{S}(1)$, 778.

787	Kosaka, Y., & Xie, SP. (2013). Recent global-warming hiatus tied to equatorial Pa-
788	cific surface cooling. <i>Nature</i> , 501(7467), 403–407.
789	Lehtinen, J., Munkberg, J., Hasselgren, J., Laine, S., Karras, T., Aittala, M., & Aila,
790	T. (2018). Noise2Noise: Learning image restoration without clean data. $arXiv$
791	preprint arXiv:1803.04189.
792	Lenssen, N. J., Schmidt, G. A., Hansen, J. E., Menne, M. J., Persin, A., Ruedy,
793	R., & Zyss, D. (2019). Improvements in the GISTEMP uncertainty model.
794	Journal of Geophysical Research: Atmospheres, 124(12), 6307–6326.
795	Li, Yu, Y., Tang, Y., Lin, P., Xie, J., Song, M., others (2020). The flexible global
796	ocean-atmosphere-land system model grid-point version 3 (FGOALS-g3): de-
797	scription and evaluation. Journal of Advances in Modeling Earth Systems,
798	12(9), e2019MS002012.
799	Li, S., & Huang, P. (2022). An exponential-interval sampling method for evaluat-
800	ing equilibrium climate sensitivity via reducing internal variability noise. Geo -
801	science Letters, $9(1)$, 1–10.
802	Maher, N., Milinski, S., Suarez-Gutierrez, L., Botzet, M., Dobrynin, M., Kornblueh,
803	L., others (2019) . The Max Planck Institute Grand Ensemble: enabling
804	the exploration of climate system variability. Journal of Advances in Modeling
805	Earth Systems, 11(7), 2050–2069.
806	Marini, C., & Frankignoul, C. (2014). An attempt to deconstruct the Atlantic multi-
807	decadal oscillation. Climate dynamics, 43, 607–625.
808	Masson-Delmotte, V., Zhai, P., Pirani, A., Connors, S. L., Péan, C., Berger, S.,
809	others (2021). Climate change 2021: the physical science basis. Contribution of
810	working group I to the sixth assessment report of the intergovernmental panel
811	on climate change, 2.
812	Meehl, G. A., Hu, A., Arblaster, J. M., Fasullo, J., & Trenberth, K. E. (2013).
813	Externally forced and internally generated decadal climate variability associ-
814	ated with the Interdecadal Pacific Oscillation. $Journal of Climate, 26(18),$
815	7298–7310.
816	Menary, M. B., Robson, J., Allan, R. P., Booth, B. B., Cassou, C., Gastineau, G.,
817	\ldots others (2020). Aerosol-forced AMOC changes in CMIP6 historical simula-
818	tions. Geophysical Research Letters, 47(14), e2020GL088166.
819	Menary, M. B., & Wood, R. A. (2018). An anatomy of the projected north atlantic

-34-

820	warming hole in cmip5 models. Climate Dynamics, $50(7-8)$, $3063-3080$.
821	Neelin, J. D., Battisti, D. S., Hirst, A. C., Jin, FF., Wakata, Y., Yamagata, T., &
822	Zebiak, S. E. (1998). ENSO theory. Journal of Geophysical Research: Oceans,
823	103(C7), 14261-14290.
824	Newman, M., Alexander, M. A., Ault, T. R., Cobb, K. M., Deser, C., Di Lorenzo,
825	E., others (2016). The Pacific decadal oscillation, revisited. Journal of
826	Climate, 29(12), 4399-4427.
827	O'Shea, K., & Nash, R. (2015). An introduction to convolutional neural networks.
828	arXiv preprint arXiv:1511.08458.
829	Parker, D., Folland, C., Scaife, A., Knight, J., Colman, A., Baines, P., & Dong, B.
830	(2007). Decadal to multidecadal variability and the climate change back-
831	ground. Journal of Geophysical Research: Atmospheres, 112(D18).
832	Parsons, L. A., Brennan, M. K., Wills, R. C., & Proistosescu, C. (2020). Magnitudes
833	and spatial patterns of interdecadal temperature variability in CMIP6. Geo-
834	$physical\ Research\ Letters,\ 47(7),\ e2019 GL086588.$
835	Rodgers, K. B., Lin, J., & Frölicher, T. L. (2015). Emergence of multiple ocean
836	ecosystem drivers in a large ensemble suite with an Earth system model. Bio -
837	$geosciences, \ 12(11), \ 3301-3320.$
838	Ronneberger, O., Fischer, P., & Brox, T. (2015). U-net: Convolutional networks for
839	biomedical image segmentation. In Medical image computing and computer-
840	$assisted\ intervention-miccai\ 2015:\ 18th\ international\ conference,\ munich,$
841	germany, october 5-9, 2015, proceedings, part iii 18 (pp. 234–241).
842	Schmidt, A., Mills, M. J., Ghan, S., Gregory, J. M., Allan, R. P., Andrews, T.,
843	others (2018). Volcanic radiative forcing from 1979 to 2015. Journal of
844	Geophysical Research: Atmospheres, 123(22), 12491–12508.
845	Sherwood, S., Webb, M. J., Annan, J. D., Armour, K. C., Forster, P. M., Har-
846	greaves, J. C., others (2020). An assessment of Earth's climate sen-
847	sitivity using multiple lines of evidence. $Reviews of Geophysics, 58(4),$
848	e2019RG000678.
849	Smith, C. J., & Forster, P. M. (2021). Suppressed late-20th century warming in
850	CMIP6 models explained by forcing and feedbacks. Geophysical Research Let-
851	ters, 48(19), e2021GL094948.

Smith, C. J., Kramer, R. J., Myhre, G., Alterskjær, K., Collins, W., Sima, A., ...
853	others (2020). Effective radiative forcing and adjustments in CMIP6 models.
854	Atmospheric Chemistry and Physics, 20(16), 9591–9618.
855	Smith, D. M., Gillett, N. P., Simpson, I. R., Athanasiadis, P. J., Baehr, J., Bethke,
856	I., \ldots others (2022). Attribution of multi-annual to decadal changes in the
857	climate system: The Large Ensemble Single Forcing Model Intercomparison
858	Project (LESFMIP). Frontiers in Climate, 4.
859	Solomon, A., Goddard, L., Kumar, A., Carton, J., Deser, C., Fukumori, I., oth-
860	ers (2011). Distinguishing the roles of natural and anthropogenically forced
861	decadal climate variability: implications for prediction. Bulletin of the Ameri-
862	can Meteorological Society, 92(2), 141–156.
863	Steinman, B. A., Mann, M. E., & Miller, S. K. (2015). Atlantic and Pacific mul-
864	tidecadal oscillations and Northern Hemisphere temperatures. Science,
865	347(6225), 988-991.
866	Sun, L., Alexander, M., & Deser, C. (2018). Evolution of the global coupled climate
867	response to Arctic sea ice loss during 1990–2090 and its contribution to climate
868	change. Journal of Climate, 31(19), 7823–7843.
869	Swart, N. C., Fyfe, J. C., Hawkins, E., Kay, J. E., & Jahn, A. (2015). Influence of
870	internal variability on Arctic sea-ice trends. Nature Climate Change, $5(2)$, 86–
871	89.
872	Swingedouw, D., Bily, A., Esquerdo, C., Borchert, L. F., Sgubin, G., Mignot, J., &
873	Menary, M. (2021). On the risk of abrupt changes in the north atlantic sub-
874	polar gyre in cmip6 models. Annals of the New York Academy of Sciences,
875	1504(1), 187-201.
876	Taylor, K. E., Stouffer, R. J., & Meehl, G. A. (2012). An overview of CMIP5 and
877	the experiment design. Bulletin of the American meteorological Society, $93(4)$,
878	485–498.
879	Tebaldi, C., Debeire, K., Eyring, V., Fischer, E., Fyfe, J., Friedlingstein, P., oth-
880	ers (2020). Climate model projections from the scenario model intercomparison
881	project (ScenarioMIP) of CMIP6. Earth System Dynamics Discussions, 2020,
882	1-50.
883	Tian, C., Fei, L., Zheng, W., Xu, Y., Zuo, W., & Lin, CW. (2020). Deep learning
884	on image denoising: An overview. Neural Networks, 131, 251–275.
885	Ting, M., Kushnir, Y., Seager, R., & Li, C. (2009). Forced and internal twentieth-

-36-

886	century SST trends in the North Atlantic. J. Climate, 22, 1469–1481.
887	Van Vuuren, D. P., Edmonds, J., Kainuma, M., Riahi, K., Thomson, A., Hibbard,
888	K., others (2011). The representative concentration pathways: an overview.
889	Climatic change, 109, 5–31.
890	Vincent, L., Zhang, X., Brown, R., Feng, Y., Mekis, E., Milewska, E., Wang, X.
891	(2015). Observed trends in Canada's climate and influence of low-frequency
892	variability modes. Journal of Climate, 28(11), 4545–4560.
893	Wang, C., & Picaut, J. (2004). Understanding ENSO physics—A review. Earth's
894	Climate: The Ocean–Atmosphere Interaction, Geophys. Monogr, 147, 21–48.
895	Wills, R. C., Battisti, D. S., Armour, K. C., Schneider, T., & Deser, C. (2020). Pat-
896	tern recognition methods to separate forced responses from internal variability
897	in climate model ensembles and observations. Journal of Climate, $33(20)$,
898	8693–8719.
899	Yamashita, R., Nishio, M., Do, R. K. G., & Togashi, K. (2018). Convolutional neu-
900	ral networks: an overview and application in radiology. Insights into imaging,
901	9,611-629.
902	Zelinka, M. D., Myers, T. A., McCoy, D. T., Po-Chedley, S., Caldwell, P. M., Ceppi,
903	P., Taylor, K. E. (2020). Causes of higher climate sensitivity in CMIP6
904	models. Geophysical Research Letters, $47(1)$, e2019GL085782.
905	Zelinka, M. D., Zhou, C., & Klein, S. A. (2016). Insights from a refined decomposi-
906	tion of cloud feedbacks. Geophysical Research Letters, $43(17)$, $9259-9269$.
907	Zhang, R. (2007). Anticorrelated multidecadal variations between surface and sub-
908	surface tropical north atlantic. Geophysical Research Letters, $34(12)$.
909	Zhang, R., Sutton, R., Danabasoglu, G., Kwon, YO., Marsh, R., Yeager, S. G.,
910	Little, C. M. (2019). A review of the role of the Atlantic meridional over-
911	turning circulation in Atlantic multidecadal variability and associated climate
912	impacts. Reviews of Geophysics, 57(2), 316–375.

Separation of internal and forced variability of climate using a U-Net

Constantin Bône¹², Guillaume Gastineau¹, Sylvie Thiria¹, Patrick Gallinari²³and Carlos Mejia¹

5	$^1 \mathrm{UMR}$ LOCEAN, IPSL, Sorbonne Université, IRD, CNRS, MNHN
6	$^2 \mathrm{UMR}$ ISIR, Sorbonne Université, CNRS, INSERM
7	³ Criteo AI Lab

Key Points:

1

2

3

4

We present a new method to separate the forced and internal variability of the surface air temperature. We utilise a U-Net trained with global climate models outputs and implement a noise to noise methodology to eliminate internal variability. The results are assessed through the utilisation of very large ensemble simulations of two distinct climate models.

 $Corresponding \ author: \ Constantin \ B\hat{o}ne, \ {\tt constantin.bone@sorbonne-universite.fr}$

15 Abstract

The internal variability pertains to fluctuations originating from processes inherent to 16 the climate component and their mutual interactions. On the other hand, forced vari-17 ability delineates the influence of external boundary conditions on the physical climate 18 system. A methodology is formulated to distinguish between internal and forced vari-19 ability within the surface air temperature. The noise-to-noise approach is employed for 20 training a neural network, drawing an analogy between internal variability and image 21 noise. A large training dataset is compiled using surface air temperature data spanning 22 from 1901 to 2020, obtained from an ensemble of Atmosphere-Ocean General Circula-23 tion Model (AOGCM) simulations. The neural network utilized for training is a U-Net, 24 a widely adopted convolutional network primarily designed for image segmentation. To 25 assess performance, comparisons are made between outputs from two single-model initial-26 condition large ensembles (SMILEs), the ensemble mean, and the U-Net's predictions. 27 The U-Net reduces internal variability by a factor of four, although notable discrepan-28 cies are observed at the regional scale. While demonstrating effective filtering of the El 29 Niño Southern Oscillation, the U-Net encounters challenges in areas dominated by forced 30 variability, such as the Arctic sea ice retreat region. This methodology holds potential 31 for extension to other physical variables, facilitating insights into the enduring changes 32 triggered by external forcings over the long term. 33

34

Plain Language Summary

To comprehensively grasp future climate change, it becomes imperative to differ-35 entiate between forced variability and internal climate variability. Internal variability refers 36 to the climate's variations driven by the chaotic nature of geophysical fluids. Conversely, 37 forced variability denotes changes prompted by external forcings, predominantly alter-38 ations in radiative forcing, primarily due to anthropogenic activities. Here, a novel ap-39 proach is introduced for filtering internal variability through the utilisation of a convo-40 lutional neural network. This neural network is trained using a noise-to-noise method-41 ology, targeting the filtration of internal variability from surface air temperature outputs 42 of climate models or observational data. Internal variability is treated analogously to noise 43 within an image, which is removed to restore the "true image," corresponding to forced 44 variability in our case. This method capitalises on the data generated by state-of-the-45 art climate models through the coupled model intercomparison project (CMIP). To val-46

-2-

idate this methodology, we assess its performance using very large ensembles of climate
model simulations, enabling precise estimation of forced variability. Our findings demonstrate a reduction in internal variability by a factor of four, accompanied by notable regional variations.

51 **1** Introduction

The phenomenon of climate warming is characterized by an elevated surface air tem-52 perature, notably reaching a pivotal juncture during the latter half of the twentieth cen-53 tury (Eyring et al., 2021). Nevertheless, the observed anomalies in surface air temper-54 ature arise from a dual spectrum of variabilities. The first source of variability is due to 55 the effect of the external forcings, such as the increase in the greenhouse gases concen-56 tration, the variations of concentration in anthropogenic and natural aerosols, the fluc-57 tuations in solar variability or volcanic eruptions and the land-use changes. The related 58 variability is designated as the forced variability. The second source of variability is com-59 ing from processes internal to the atmosphere, oceans, cryosphere and land or the inter-60 actions between them (Cassou et al., 2018). Subsequently, this form of variability is re-61 ferred to as 'internal variability,' encapsulating its inception within the climate system 62 and its persistence even without alterations in external forcings. Despite the overarch-63 ing dominance of forced variability in shaping the broad-scale and long-term trajectory 64 of surface air temperature across the 1900-2020 timeframe (Deser et al., 2012; Kay et 65 al., 2015), a comprehensive understanding of the distinct contributions of internal and 66 forced variability remains elusive. Internal variability takes center stage in briefer tem-67 poral scales and smaller spatial dimensions. For instance, the leading mode of internal 68 variability in global air surface temperature manifests as the El Niño Southern Oscilla-69 tion (ENSO), characterized by significant anomalies in the equatorial Pacific Ocean, ac-70 companied by distant teleconnections, and a prevailing cycle spanning two to seven years 71 (Wang & Picaut, 2004). Additionally, the interdecadal Pacific variability (Newman et 72 al., 2016) and the Atlantic Multidecadal variability (Zhang et al., 2019) wield the capac-73 ity to influence climate dynamics across the decadal to multidecadal spectrum. A no-74 table example involves the deceleration in the global warming rate experienced during 75 2002-2012, commonly referred to as the global warming hiatus, which has been robustly 76 linked to Interdecadal Pacific Variability (Meehl et al., 2013; Kosaka & Xie, 2013; Eng-77 land et al., 2014). Lastly, internal variability exercises influence even over centennial and 78

-3-

multi-centennial spans (Jiang et al., 2021; S. Li & Huang, 2022) exerting substantial impact on trends within the 1900-2015 interval (Bonnet et al., 2022).

The distinction between forced variability and internal variability is essential for conducting detection and attribution studies, enabling accurate estimation and simulation of the climate's reaction to alterations in radiative forcing. Moreover, this differentiation aids in recognizing and comprehending internal climate variability. Nevertheless, the availability of instrumental observations is limited to the period since 1850, and the relatively brief duration of these observations presents challenges in effectively and confidently discerning internal variability.

For identifying both internal and forced variability, linear trends (Swart et al., 2015; 88 Vincent et al., 2015) or quadratic trends (Enfield & Cid-Serrano, 2010) have been em-89 ployed to characterize forced variability. However, linear or quadratic trends inadequately 90 capture the temporal evolution of temperature, particularly failing to account for the abrupt 91 cooling subsequent to significant volcanic eruptions, which hold significant climate im-92 pact (Schmidt et al., 2018). Additional approaches include the application of Empiri-93 cal Orthogonal Functions (EOF) analysis (Parker et al., 2007), low-frequency pattern 94 filtering (Wills et al., 2020), and linear inverse models (Marini & Frankignoul, 2014). These 95 techniques deconstruct forced variability into a combination of modes featuring distinct 96 patterns and corresponding time series. Regression analysis of the global mean surface 97 temperature (GMST) has also been employed, although this may inadvertently estab-98 lish misleading links between the Atlantic and Pacific basins (Frankignoul et al., 2017; 99 Deser & Phillips, 2023). However, a comprehensive and systematic examination of these 100 methodologies remains notably absent. 101

Climate model simulations have been employed to overcome the limitations of sparse 102 observation sampling. Conducting an ensemble of climate model simulations with diverse 103 initial conditions enables estimation of forced variability via the ensemble mean. This 104 approach effectively mitigates the variance linked to internal variability by a factor of 105 n, where n signifies the ensemble's size (Harzallah & Sadourny, 1995; Hawkins & Sut-106 ton, 2009; Ting et al., 2009; Solomon et al., 2011; Deser et al., 2014; Frankcombe et al., 107 2015). As a result, modeling centers have undertaken substantial ensembles with over 108 20 or 30 ensemble members (Jeffrey et al., 2013; Rodgers et al., 2015; Sun et al., 2018; 109 Deser et al., 2020). These large ensembles are commonly referred to as Single-Model Initial-110

-4-

Condition Large Ensembles (SMILE; Deser et al. (2020)). Multiple SMILE initiatives 111 have been undertaken using models such as CCSM3 (Collins et al., 2006), CCSM4 (Gent 112 et al., 2011), CESM (Kay et al., 2015), MPI-ESM (Maher et al., 2019), FGOALS-g3 (Li 113 et al., 2020), CanESM2 (Chylek et al., 2011), and IPSL-CM6A-LR (Bonnet et al., 2021), 114 among others. This offers a valuable dataset for crafting methodologies dedicated to the 115 disentanglement of forced and internal variability. Notably, employing members of a large 116 ensemble model as surrogate observations allows for a comparison of results with the en-117 semble mean. Differences primarily mirror residual internal variability or limitations in-118 herent in the method. 119

Nevertheless, the forced variability estimated through an ensemble mean remains 120 contingent upon the specific climate model employed. These climate models carry sub-121 stantial uncertainties, particularly in terms of their climate sensitivity (Sherwood et al., 122 2020), often attributed to factors like uncertain cloud retroaction which significantly im-123 pact equilibrium climate sensitivity (Zelinka et al., 2016). Additionally, significant un-124 certainties surround historical emissions and the linked radiative forcing from aerosols 125 (Menary et al., 2020; C. J. Smith & Forster, 2021). Moreover, the internal variability ex-126 hibited by different models also varies significantly (Parsons et al., 2020). 127

Several methodologies have been devised to harness data from diverse climate mod-128 els, as employing a multi-model approach holds the potential to alleviate the uncertain-129 ties inherent in individual climate models. Multi-model ensemble means are widely adopted 130 for estimating the forced signal (Steinman et al., 2015). Notably, techniques such as the 131 signal-to-noise-maximizing empirical orthogonal functions (Ting et al., 2009; Wills et al., 132 2020) and the discriminant analysis and maximization of the average predictability time 133 (DelSole et al., 2011) have been put forth to extract forced variability with superior ef-134 ficacy compared to ensemble means. Furthermore, scaling techniques that adjusts the 135 forced signal from models using observational data have been proposed. Among these 136 methodologies are fingerprinting methods grounded in linear regression, commonly ap-137 plied for detecting and attributing climate change with a unified forcing that encapsu-138 lates the influence of all external forcings (Hasselmann, 1993; Allen & Tett, 1999; Allen 139 & Stott, 2003). More recently, the use of scaling factors was also proposed by Frankcombe 140 et al. (2015). 141

-5-

This paper introduces an alternative approach to distinguishing internal and forced 142 variability using climate model data, employing a non-linear method that takes into ac-143 count the spatio-temporal data covariances. This method is rooted in a neural network 144 trained on data from Atmosphere-Ocean General Circulation Models (AOGCMs). Among 145 the areas where neural networks have excelled is image analysis (Egmont-Petersen et al., 146 2002). One of the prominent applications of neural networks in image processing is im-147 age denoising, involving the elimination of noise from an image to restore its true form 148 (Ilesanmi & Ilesanmi, 2021; Tian et al., 2020). In this context, internal variability is treated 149 as noise. It is demonstrated that machine learning image denoising methodologies can 150 subsequently isolate forced variability. The internal variability is eliminated, leaving be-151 hind a quantifiable residue. This method leverages the temporal and spatial information 152 inherent in climate models to establish the weights and biases of a neural network. With 153 these parameters in place, the neural network is also employed with observations to delve 154 into and attribute the progression of climate change since 1905 to 2016. To the best of 155 our knowledge, this represents the pioneering application of a dedicated neural network 156 for the purpose of disentangling internal and forced variability. 157

The structure of this paper is as follows: Section 2 outlines the data utilized. Section 3 introduces the method anchored in a neural network. Section 4 assesses the method's performance. In Section 5, the neural network method is applied to observations. Lastly, Section 6 offers the conclusion and discussion.

162 **2 Data**

¹⁶³ 2.1 Observations

The gridded monthly Surface Air Temperature anomaly (SAT) from 1901 to 2020, as provided by GISS Surface Temperature Analysis version 4 (GISTEMP; Hansen et al. (2010); Lenssen et al. (2019)), is employed in this study. GISTEMP amalgamates meteorological station data over land (NOAA GHCN v4) with sea surface temperature (SST) estimates from ERSST v5. This data is available on a consistent 2°x2° grid. The monthly values are aggregated to calculate annual means, and the SAT anomalies are determined using the reference period 1950-2014. 171

2.2 Climate model simulations

The monthly SAT data is sourced from historical simulations within the Coupled 172 Model Intercomparison Project Phase 5 (CMIP5; Taylor et al. (2012)) and the Coupled 173 Model Intercomparison Project Phase 6 (CMIP6; (Eyring et al., 2016)), along with sev-174 eral Single-Model Initial-Condition Large Ensembles (SMILEs) from distinct models: MPI-175 ESM (Maher et al., 2019), CSIRO-Mk3-6-0 (Collier et al., 2011), EC-Earth (Döscher et 176 al., 2021), and FGOALS-g3 (Li et al., 2020). For the historical simulations, spanning 1901 177 to 2005 (2014) for CMIP5 (CMIP6), all external forcings are integrated. These forcings 178 encompass the effects of historical greenhouse gas concentrations, anthropogenic and nat-179 ural aerosols, stratospheric ozone, solar activity, and land-use changes. Each climate model 180 delivers multiple realizations referred to as ensemble members, generated through dis-181 tinct initial conditions. From 2005 (2014 for CMIP6) until 2020, the outputs under the 182 pessimistic Representation Concentration Pathway 8.5 (RCP8.5) scenario for CMIP5 (Van Vu-183 uren et al., 2011) and the intermediate Shared Socio-economic Pathway 2 4.5 (SSP2-4.5) 184 for CMIP6 (Tebaldi et al., 2020) are employed. These simulations utilize socio-economic 185 assumptions to project future external forcing patterns. Additionally, several SMILES 186 are incorporated, employing distinct historical forcings or scenario simulations of CMIP5 187 or CMIP6 (elaborated in Table S3). While minor differences are anticipated in exter-188 nal forcing between CMIP5 and CMIP6 simulations, notable uncertainties arise in aerosol 189 emissions (C. J. Smith et al., 2020; Fyfe et al., 2021). Modest differences may also emerge 190 between the RCP8.5 (strong) and SSP2-4.5 (moderate) scenarios, particularly until 2020, 191 where actual forcings mirror observed forcings to a considerable extent (Masson-Delmotte 192 et al., 2021). 193

The count of members accessible for scenario simulations is fewer compared to the 194 historical counterparts. Therefore, we extended the outputs from historical experiments 195 using the scenario ensemble member of the same model with the same number identi-196 fication. In case the number identification is lacking, we select randomly an scenario en-197 semble member of the same climate model. 198

199

All monthly data are aggregated into annual means. Subsequently, the SAT anomalies are computed for each ensemble member using 1950-2014 as a reference period. This 200 furnishes a multi-model ensemble comprising 801 members derived from 47 AOGCMs. 201 Subsequently, the concatenated historical and scenario members are harnessed within 202

-7-

the 1901-2020 timeframe. All model data is regridded using bilinear interpolation on the horizontal grid from GISTEMP. The details pertaining to the climate model names, ensemble sizes, and the names of the employed scenario simulations are elucidated in Tabs. S1, S2, and S3.

207

2.3 Validation of the data set

The forced variability simulated within the multi-model ensemble is succinctly ex-208 amined for two specific data subsets. We investigate the MPI-ESM and FGOALS-g3 cli-209 mate models from SMILE, as they have a very large size of 100 and 115 members, re-210 spectively, which largely exceed the size of other model ensembles. Anticipatedly, the es-211 timated forced variability derived from the ensemble mean for each of these models is 212 expected to be accurate, as the reduction in variance attributed to internal variability 213 reaches 100 and 115, respectively. For instance, Deser et al. (2012, 2014) demonstrated 214 that identifying regional climate responses on time scales of several decades may neces-215 sitate between 10 to 40 members. Specifically, to detect a change in SAT between the 216 decades 2005-2014 and 2028-2037 on a global scale, the use of 3 to 6 members is requi-217 site. This requirement can surge beyond 10 for local analyses such as in North Amer-218 ica. Subsequently, the data originating from these two models is subsequently employed 219 to appraise the outcomes of the neural network model in section 4.1. 220

We utilize the ensemble mean to characterize the forced variability and employ the 221 standard deviations from the ensemble members for evaluating the internal variability. 222 Figure 1 illustrates the standard deviation of the SAT deviation from the ensemble mean 223 for FGOALS-g3 and MPI-ESM. The variability in SAT is more pronounced over land 224 surfaces ($\sim 0.3^{\circ}$ C) compared to oceans ($\sim 0.1^{\circ}$ C), consistent with the lower thermal in-225 ertia of land. Notably, substantial variability (ranging from approximately 1.5°C to 2.5°C) 226 is observed over regions coinciding with the sea ice edge, such as the Bering Sea and Nordic 227 Seas in the Northern Hemisphere, as well as the Amundsen and Weddell Seas in the South-228 ern Hemisphere. Additionally, a marked variability is observed in the equatorial Pacific 229 Ocean, with a standard deviation of 0.8°C, and this variability is more prominent in MPI-230 ESM compared to FGOALS-g3. A localized peak of variability is situated over the sub-231 polar North Atlantic, especially notable for FGOALS-g3 (reaching up to 2°C). These out-232 comes coherently reflect a significant internal variability stemming from extratropical weather 233 fluctuations over land surfaces, exhibiting local maxima around regions adjacent to the 234

-8-



Figure 1. Standard deviation of the SAT deviations from the ensemble mean for (top) MPI-ESM and (bottom) FGOALS-g3.

sea ice edge. Moreover, the variability observed in the equatorial Pacific mirrors the phenomenon of El Nino Southern Oscillation (Neelin et al., 1998).

The forced variability is estimated through the ensemble mean of each model. Sub-237 sequently, the multi-model mean (MMM) is computed by averaging the ensemble means 238 across all models, ensuring equal weight for each model. Nonetheless, MPI-ESM and FGOALS-239 g_3 are excluded from this computation, as the intention is to later compare them to the 240 MMM. To assess the prominent impact of greenhouse gas forcing, Figure 2 (a, c, e) il-241 lustrates the ensemble mean SAT anomaly for MPI-ESM, FGOALS-g3, and the MMM 242 throughout the 2010-2020 interval. Furthermore, Figure 2 (b, d, f) presents the tempo-243 ral standard deviation of the ensemble means across the period from 1901 to 2020. As 244 anticipated, all climate models project more substantial warming over land (up to 0.8°C) 245 than over oceans (approximately 0.3°C). Notably, the Arctic exhibits an amplification 246 of global warming, with warming exceeding 2°C north of 60°N. The MMM showcases an 247 average warming of 0.8°C for the 2010-2020 period, surpassing MPI-ESM (0.64°C) and 248

FGOALS-g3 (0.69°C). This aligns with the comparatively lower equilibrium climate sen-249 sitivity (ECS) of these two models (3.6°C for MPI-ESM and 2.8°C for FGOALS-g3) when 250 compared to other models employed in this study (Zelinka et al., 2020). Within the sub-251 polar Atlantic, the SAT anomalies exhibit a minimum, with negative temperatures anoma-252 lies observed in FGOALS-g3 over the Labrador Sea, or in MPI-ESM over the subpolar 253 gyre. This phenomenon, known as the North Atlantic warming hole (Keil et al., 2020), 254 is associated with a deceleration of the Atlantic meridional overturning circulation (He 255 et al., 2022). It is worth noting that such a minimum is less pronounced in the MMM, 256 presumably due to considerable uncertainties regarding the precise location of this warm-257 ing hole and the linked processes. An equivalent spatial pattern can be derived using stan-258 dard deviations, revealing values of approximately 0.3°C for the majority of global re-259 gions and higher values over land ($\sim 0.6^{\circ}$ C). Grid points located north of 60° also exhibit 260 elevated values, peaking at around 2°C in the Barents Sea for MPI-ESM or the Labrador 261 Sea for FGOALS-g3. 262

The forced variability exhibited by MPI-ESM and FGOALS-g3 diverges from that 263 of the MMM, revealing a comparatively weaker global warming trend and standard de-264 viation pattern. This divergence is particularly evident north of 60°N, where the warm-265 ing exhibits greater amplification (refer to Fig. 2), amounting to 1.54°C for MPI-ESM 266 and 1.45°C for FGOALS-g3. Local variations are also observed in regions such as the Labrador 267 Sea, Barents and Kara Sea, the Canadian archipelago, and the Bering Sea in the case 268 of FGOALS-g3. Notably, MPI-ESM similarly presents notable differences in the Barents 269 Sea. These discrepancies may arise from biases related to sea ice representation. Specif-270 ically, FGOALS-g3 depicts an excessive extent of Arctic sea ice (Li et al., 2020), which 271 in turn leads to inaccuracies in simulating the location of the sea ice edge. This discrep-272 ancy can account for spurious SAT variability attributed to the misplaced sea ice edge 273 within the Labrador Sea. The mean standard deviation of the ensemble mean registers 274 as 0.34°C for MPI-ESM and 0.43°C for FGOALS-g3, exceeding the mean standard de-275 viation of the SAT deviations of the members to the ensemble mean which is of 0.51°C 276 for MPI-ESM and 0.46°C for FGOALS-g3. This underscores that the internal variabil-277 ity is marginally more pronounced than the forced variability. 278

-10-



Figure 2. a) Ensemble mean of the air surface temperature (°C) in MPI-ESM in 2010-2020. c) Same as a) but for FGOALS-g3. e) Same as a) but for the MMM. b) Standard deviation of the ensemble mean surface air temperature (°C) in 1901-2020 for MPI-ESM. d) Same as b) but for FGOALS-g3. f) Same as b) but for the MMM.

279 3 Methods

3.1 Neural network

We design a neural network to remove the internal variability from the SAT. The 281 input data is structured with dimensions (120, 90, 180), corresponding to time spanning 282 from 1901 to 2020, latitude, and longitude, respectively. On the other hand, the output 283 holds dimensions of (112, 90, 180), encompassing the years 1905 to 2016, while maintain-284 ing the latitude and longitude dimensions intact. Notably, the output's temporal span 285 is truncated compared to the input, by excluding the initial and final four years. This 286 reduction addresses the substantial uncertainty typically observed at the dataset's end-287 points, an aspect that will be elaborated upon later. 288

A neural network's characteristics are shaped by its hyperparameters, which dictate both its architecture and training process. Our approach involves utilizing three distinct datasets, each composed of input and desired output pairs. The training dataset serves the purpose of establishing the neural network's weights and biases. Meanwhile, the validation dataset comes into play for estimating the hyperparameters. Finally, the test dataset is employed to assess the neural network's performance.

295

3.2 Constitution of the database

To construct the training dataset, we adapt a noise-to-noise methodology originally 296 introduced in Lehtinen et al. (2018). This approach was initially designed to train a neu-297 ral network in denoising images. In this method, the network is exclusively trained on 298 noisy images depicting various objects. Each object has more than one noised image de-299 picting it. In the noise to noise method, we create an input/output training database 300 that comprises pairs of noisy image combinations for identical objects. It's essential to 301 note that the network cannot effectively learn to transform a random noise realization 302 into another. Instead, the configuration is designed to approximate the mathematical 303 expectation of all noisy images associated with the same object, culminating in an es-304 timate that closely resembles the noise-free image. 305

For our application, we consider the forced spatio-temporal SAT anomalies from each climate model as distinct objects. These anomalies, inherent to each member, can be likened to noisy images, where the internal variability introduces the noise component. The ensemble members' mathematical expectation equates to the forced variabil-

ity, which can be approximated through the ensemble mean.

To create the training dataset, we follow a procedure wherein we compute pairs of 311 members for each climate model, except for MPI-ESM, FGOALS-g3, and MIROC6, which 312 are reserved for testing and validation purposes. Adopting an approach similar to Lehtinen 313 et al. (2018), we augment the dataset by introducing the ensemble mean of the climate 314 model's members as an additional member. This inclusion serves to expedite the train-315 ing process without introducing any other influences. In this process, each pair of mem-316 bers becomes an input/output pair. If we denote the number of ensemble members ob-317 tained from a specific climate model as n, this approach yields n(n+1) input/output 318 pairs per model. By accumulating such pairs from all models, the resulting training dataset 319 primarily comprises simulations characterized by the most extensive ensemble sizes (namely 320 IPSL-CM6A-LR, CanESM5, CNRM-CM6-1, and ACCESS-ESM1-5). 321

To create the validation set, we employ the ensemble simulation data from the MIROC6 model, which ranks as the third-largest ensemble in terms of size (with n = 50 members). For this purpose, we designate the ensemble members as inputs, while the ensemble mean spanning the period from 1905 to 2016 serves as the desired output.

To form the test dataset, we draw upon data derived from the FGOALS-g3 and MPI-ESM models, leveraging their extensive ensemble sizes of n = 110 and n = 100respectively. Subsequently, we proceed to make comparisons between the outputs of the neural network obtained from ensemble members and their corresponding ensemble means for both of these models.

The conclusions drawn from these tests and validation processes may exhibit some dependence on the specific model being analyzed, as alternative models could yield varying outcomes. Nevertheless, this approach has been chosen due to its simplicity and its potential to mitigate the impact of any remaining internal variability.

335 **3.3 U-Net**

Convolutional neural networks (CNNs, Yamashita et al. (2018)) constitute a category of non-linear neural networks, notably applied in tasks related to imagery (O'Shea & Nash, 2015). A distinctive attribute of CNNs is their utilization of convolutional layers, which incorporate a trainable kernel that slides across the input data.

-13-



Figure 3. Schematic of the U-Net. The arrows represent the operations within the network. The numbers shows the dimension of the data and the number of filters used.

In this context, a U-Net architecture is employed, which falls within the realm of 340 CNNs. Originally introduced by Ronneberger et al. (2015) for image segmentation, the 341 U-Net structure has gained widespread popularity in image-related analyses such as de-342 noising (Ilesanmi & Ilesanmi, 2021; Tian et al., 2020). The U-Net architecture is char-343 acterized by its inclusion of a contracting path and an expansive path, which collectively 344 give rise to its characteristic U shape (refer to Fig. 3). The contracting path adheres to 345 a conventional design of a convolutional network, featuring numerous convolutional lay-346 ers, each followed by an activation function and a max-pooling operation. As the con-347 tracting path advances, spatial information is diminished while feature information is 348 enriched. Conversely, the expansive path amalgamates feature and spatial information 349 through a sequence of up-convolutions and concatenations with high-resolution features 350 derived from the contracting path. 351

The U-Net architecture employed in this study shares similarities with the design proposed by Ronneberger et al. (2015). However, a modification is made by replacing the 2-dimensional convolutional layers with 3-dimensional counterparts. This alteration is introduced to encompass not only the spatial dimension but also the temporal dimension of the data. The selected activation function is the hyperbolic tangent. Additionally, adaptations have been made to the output layer to accommodate an output com-

-14-

prising 112 time steps. The neural network is comprised of a total of 5,659,009 trainable
 parameters.

A batch size of 8 is chosen, and the optimization process employs the Adam op-360 timizer with a learning rate of 0.001. To ensure proper application of the CNN to the 361 data, padding is introduced. This involves extending the image by appending zero val-362 ues at its edges. For the longitudinal dimension, which is periodic, the zero padding only 363 results in a slight discontinuity at 180°E, the edge of the data. Indeed, due to the na-364 ture of convolutional layers, a U-Net has more difficulty processing information located 365 at the edge of the data. This is the reason why we excluded the initial and final four years 366 (1901-1904 and 2017-2020) in the U-Net's outputs. The chosen cost function is the root 367 mean squared error (RSME), calculated using an area-weighted mean of the gridded data. 368

The validation dataset is utilized to determine the optimal values for two key hy-369 perparameters: the number of epochs and the number of filters used in the convolutional 370 layers. The term "number of filters" pertains to the thickness of the convolutional lay-371 ers. The number of epochs refers to how many times the training dataset is processed 372 during the training phase. These hyperparameters are selected to minimize the root mean 373 squared error (RMSE) using the validation dataset. Examination of the validation RMSE 374 for different values of epochs and layer thickness reveals a consistent pattern (see Fig. 375 S1): a significant reduction in RMSE occurs in the initial epochs, followed by a grad-376 ual increase. As a result, we settle on a layer thickness of 16 for the first layer (as shown 377 in Fig. 3) and a total of 32 epochs. 378

379 **3.4 Example**

Figure 4 provides an illustrative example featuring two randomly selected ensem-380 ble members from MPI-ESM and FGOALS-g3. The comparison focuses on the SAT at 381 the year 2016, depicted in the top panels, as well as the resulting output generated by 382 the neural network in 2016 (centre panels), juxtaposed against the ensemble mean anomaly 383 for the same year (bottom panels). The anticipated impact of elevated greenhouse gas 384 concentrations in 2016 is evident in the SAT of both MPI-ESM and FGOALS-g3 mem-385 bers, which exhibit warm anomalies. However, the internal variability introduces anoma-386 lies that surpass those of the ensemble mean in numerous regions, accompanied by some 387 negative anomalies in other areas. To elaborate, an instance of cooling is simulated across 388

-15-

the Equatorial Pacific Ocean, possibly linked to a La Niña event in the case of MPI-ESM. 389 The same ensemble member displays cooling over land in equatorial Africa, South-Eastern 390 Asia, and Australia, as well as in extratropical zones like the North Atlantic Ocean and 391 the Weddell Sea. In the example from FGOALS-g3, cold anomalies emerge over the Nordic 392 Seas and the Labrador Sea. Such cooling diverges from the ensemble average, which ex-393 hibits a relatively uniform warming pattern across the globe, with a more pronounced 394 effect over landmasses. Notably, the Arctic and its environs experience heightened warm-395 ing compared to other global regions, due to polar amplification. Conversely, minimal 396 warming is observed in the Southern Ocean and the subpolar North Atlantic Ocean, and 397 even a cooling tendency is noted in the Northern Atlantic warming hole. 398

The SAT obtained from the U-Net's output, utilizing the same ensemble member 399 as input, exhibits a pattern strikingly similar to that of the ensemble mean (compare cen-400 tre and bottom panels). In both instances, the pattern is relatively uniform, albeit with 401 heightened warming observed over land areas, coupled with an Arctic Amplification phe-402 nomenon. This suggests that the internal variability—such as the influence of ENSO events 403 or the effects of prolonged weather patterns over continents—has been successfully elim-404 inated. The regions displaying subdued warming or cooling tendencies are replicated, 405 although the exact positioning and intensity might not precisely match those of the en-406 semble mean in certain areas, particularly the Southern and subpolar North Atlantic. 407 It's worth noting a minor discontinuity at 180°E resulting from the padding process. 408

The performance of the method is quantified more systematically in the next sec-tion.

411

4 The U-Net as an internal variability filter

The U-Net was applied to every member of FGOALS-g3 and MPI-ESM. We then compare the results obtained with the respective ensemble mean of these two climate models.

Figures 5a and 5b illustrate the root mean squared error (RMSE) between the outcomes generated by the U-Net and the corresponding ensemble mean for the time period of 1905-2016. Notably, the discrepancies in U-Net's predictions are not uniformly distributed across space. The RMSE values fall within the range of 0.05°C to 0.5°C. The discrepancies generally remain below 0.2°C in tropical regions, except for instances over

-16-



Figure 4. (First column) Anomalies of SAT in a randomly chosen member of MPI-ESM, the associated U-Net output and ensemble mean in 2016. (Second column) Same as the first column but for a randomly chosen ensemble member for FGOALS-g3.

Western Africa in the MPI-ESM model. In contrast, the largest errors are concentrated 420 in polar areas, encompassing the Nordic Seas, Labrador Sea, and Bering Sea. Moreover, 421 sizable errors are also evident over the Southern Ocean and the continents of the North-422 ern Hemisphere situated above 45°N. These high-error regions correspond to locales char-423 acterized by substantial internal variability (refer to Figure 1). Nevertheless, it is note-424 worthy that the errors produced by the U-Net are approximately five times smaller than 425 the actual internal variability. Between the years 1996 and 2016, both ensemble results 426 exhibit a warming trend that is roughly 0.1°C lower in the U-Net results when compared 427 to the ensemble mean (as observed in Figs. 5cd). This difference is indicated by the nearly 428 consistent negative divergence situated between latitudes 45°N and 45°S. 429

The prevailing trend of systematic underestimation is, however, disrupted by an 430 exception involving the subpolar Atlantic and the Southern Ocean, where an overesti-431 mation of warming is observed. This overestimation is particularly conspicuous in the 432 FGOALS-g3 model, with warming anomalies extending to approximately 1°C over the 433 Labrador Sea and 0.5°C over the Bering Sea. This divergence from the ensemble mean 434 highlights the limited capacity of the neural network to accurately predict forced changes 435 within the subpolar North Atlantic, which is a region that exhibits inconsistent surface 436 temperature shifts across models (Swingedouw et al., 2021). The neural network's per-437 formance is restricted due to this discrepancy among models, which hampers its abil-438 ity to discern the specific features of each climate model. For example, in the case of FGOALS-439 g3, the extensive anomalies in the Labrador and Bering Seas are not mirrored in the multi-440 model mean (see Figure 2). It's also plausible that the substantial internal variability 441 observed in these regions poses a challenge for accurate removal by the neural network 442 (refer to Figure 1). This underestimation extends to the continents, with a greater im-443 pact on South America, Africa, and Australia in the tropics, as well as North America 444 and Northern Siberia in boreal regions. The degree of underestimation reaches 0.15°C 445 for MPI-ESM and 0.13°C for FGOALS-g3 in these regions. 446

Figures 6c and 6d illustrate the temporal evolution of the global surface air temperature (GSAT) for both the MPI-ESM and FGOALS-g3 models, before and after applying the U-Net correction. The range of data variability is portrayed by a 90% confidence interval assuming an Gaussian distribution. The forced variability's temporal trend extracted via ensemble mean (depicted by the red line) is effectively captured by the U-Net outputs (represented by the blue line and blue shading).

-18-



Figure 5. a) Root mean square difference of the surface air temperature, in C°, between the outputs of the U-Net and the mean ensemble in MPI-ESM, calculated across the members and all years in 1905-2016 b. b) Same as a) but for FGOALS-g3 c) Difference of the time mean SAT anomaly during 1996-2016, in °C, between the mean output of the U-Net and the corresponding ensemble mean, for MPI-ESM. d) Same as c) but for FGOALS-g3

- From 1905 to 2016, a GSAT rise is observed, aligning with the anticipated shifts in radiative forcing (Gulev et al., 2021). Additionally, a cooling pattern emerges a few years subsequent to the significant volcanic eruptions of Agung (1963), El Chichón (1982), and Pinatubo (1991), a phenomenon accurately estimated by the U-Net. This outcome aligns with expectations based on climate models incorporating volcanic aerosol emissions. Impressively, the U-Net's outputs exhibit a marginal spread, reduced approximately tenfold, indicating a substantial removal of internal variability.
- Nonetheless, the U-Net results exhibit anomalies with a slightly diminished am-460 plitude compared to the ensemble mean. The spread of the U-Net outputs is also ap-461 proximately twice as wide at the time series' beginning and end. The distribution of spa-462 tially averaged RMSE values within 90°S-90°N, comparing all U-Net outputs to the en-463 semble mean (depicted in Fig. 6a and 6b as blue histograms), reveals errors of around 464 0.12°C in MPI-ESM and 0.13°C in FGOALS-g3. Additionally, we examine the RMSE 465 values when averaging within 60°N-90°N, as Fig. 5ab suggests that errors are most pro-466 nounced in this region (illustrated in Fig. 6ab as red histograms). Errors north of 60°N 467

are approximately twice as substantial as global averages, with an average error of around 468 0.23°C in MPI-ESM and 0.26°C in FGOALS-g3. In Fig. 6ef, the internal variability ob-469 served when averaging the SAT north of $60^{\circ}N$ (as depicted by the red shading) is con-470 siderable in the raw model outputs (around 0.8° C). The ensemble mean SAT anomalies 471 in this region increase from approximately -1°C in the early twentieth century to about 472 1.2°C in 2010. The temporal evolution of the SAT north of 60°N demonstrates notable 473 similarity between the ensemble mean and the ensemble mean of U-Net outputs, with 474 a roughly 10-fold reduction in spread. However, the amplitude of the anomalies is slightly 475 underestimated, with a reduction of around 0.3°C in negative anomalies in the U-Net 476 output between 1905 and 1930 in MPI-ESM. For FGOALS-g3, the SAT is underestimated 477 by around $0.2^\circ\mathrm{C}$ during 1970-1990. 478

In Figure S2, the quadratic errors between the mean ensemble members and the 479 U-Net output are presented for each year, with global (90°S-90°N) and north of 60°N av-480 erages considered for both MPI-ESM and FGOALS-g3. Notably, the RMSE exhibits el-481 evated values during the initial and final years, characterized by peaks around the years 482 1975-1985 in both models. This pattern underscores the presence of substantial uncer-483 tainties at the data's onset and conclusion. When applying the 1900-2020 period for the 484 output (without excluding the first and last four years), the errors actually surpass those 485 portrayed in Figure S2, a fact that elucidates the rationale for excluding the endpoints 486 in the ongoing analysis, as detailed in the methods (section 2). Moreover, the notable 487 error peak during 1975-1985 lacks a definitive explanation, although it's plausible that 488 this discrepancy could be linked to uncertainties associated with the implementation of 489 aerosol forcings, notably CMIP5 for MPI-ESM and CMIP6 for FGOALS-g3. 490

The errors exhibited by the U-Net in relation to data from FGOALS-g3 are more prominent compared to those arising from the use of MPI-ESM data. This discrepancy can be attributed to the fact that MPI-ESM's simulated forced variability aligns more closely with the training data's characteristics, on average. Specifically, the training data's forced variability is in line with that of the MMM, and MPI-ESM demonstrates a smaller root mean squared difference from the MMM compared to FGOALS-g3 (as illustrated in Fig. 2).

To assess the reduction in internal variability achieved by the U-Net, we can quantitatively measure the number of ensemble members needed to surpass the U-Net's in-

-20-



Figure 6. a) Histogram showing the distribution of the RMSE between the mean ensemble and the U-Net outputs of MPI-ESM. b) Same as a), but for FGOALS-g3. c) Time evolutions of the global mean surface air temperature, in °C, for the ensemble mean and the mean U-Net outputs for MPI-ESM. Color shade shows the spread of the time series, with 90% the ensemble members uncertainty assuming a gaussian distribution. d) Same as c) but for FGOALS-g3. e) and f) are the same as c) and d) but when averaging the SAT, in °C, north of 60°N.

dividual member results using a basic ensemble mean approach. This evaluation is con-500 ducted through a random subsampling process involving 500 sets of m members, where 501 m varies from 1 to 40, for both the FGOALS-g3 and MPI-ESM ensembles. Within each 502 subset, ensemble means are calculated. The RMSE between these subsample ensemble 503 means and the actual ensemble mean obtained from all members is then determined (de-504 picted by vertical red and blue lines in Figure 7). This RMSE computation is performed 505 across all grid points and is spatially averaged. The 90% intervals, assuming an Gaus-506 sian distribution, of the 500 subsamples are also illustrated. This analysis is done for both 507 the MPI-ESM and FGOALS-g3 ensembles across distinct geographical regions: global 508 (90°S-90°N), North Atlantic (60°W-0°E, 0°N-60°N), North Pacific (120°E-100°W, 20°N-509 60° N), Niño3 (5°N-5°S, 150°W-90°W), as well as polar regions north of 60°N and south 510 of 60°S. These chosen regions exhibit considerable forced and internal variability, as vi-511 sually demonstrated in Fig. 1 and Fig. 2. Additionally, this evaluation is extended to 512 encompass both oceanic and terrestrial areas in the 60°S-60°N band, allowing for a more 513 comprehensive understanding of the U-Net's performance. The horizontal lines in the 514 illustration correspond to the same RMSE values but for the U-Net output from each 515 individual member. The accompanying color shade represents the spread of 90% uncer-516 tainty assuming an Gaussian distribution. 517

Figure 7a visually illustrates the progression of errors within the subset of mem-518 bers as the size of the subset increases. This pattern aligns with expectations, as a larger 519 subset size leads to better estimations of forced variability and a corresponding reduc-520 tion in residual internal variability by a factor of \sqrt{n} . The distribution of U-Net outputs 521 mirrors the histograms presented in Figure 6, showing a high degree of similarity across 522 both climate models. The U-Net effectively diminishes internal variability in GSAT by 523 approximately a factor of slightly more than four, which is analogous to the residual vari-524 ability observed within subsets containing around 17 members for FGOALS-g3 and 20 525 members for MPI-ESM. When focusing on regions spanning oceans and land between 526 60°N and 60°S, the outcomes remain largely consistent, showcasing a reduction in error 527 magnitude by a factor of approximately four. This reduction corresponds closely to that 528 achieved by using a subset of 15 to 20 members. 529

The U-Net's efficacy stands out prominently over the equatorial Pacific region, as depicted in panel 7f. This region is known for being heavily influenced by the ENSO, which dominates internal variability. The U-Net achieves a substantial reduction in variabil-

-22-

ity, amounting to a factor of 5.5. This reduction is akin to the outcome of utilizing an
 ensemble mean derived from around 30 members for both MPI-ESM and FGOALS-G3.

In other regions, the variability reduction is quite similar to that found globally. 535 For instance, this consistency is observed in the North Pacific and polar regions, where 536 the required number of members for equivalent outcomes remains relatively steady. How-537 ever, in terms of removing internal variability, the U-Net showcases higher efficiency in 538 the context of MPI-ESM for most scenarios. This pattern holds true except for the North 539 Atlantic, where a notable deviation is observed: a set of 15 members is necessary in MPI-540 ESM to achieve results equivalent to the U-Net (~4-fold reduction in residual variabil-541 ity), while merely 5 members suffice for FGOALS-g3 (halving of the residual variabil-542 ity). 543

The variation in performance between FGOALS-g3 and MPI-ESM might arise from dissimilarities in their internal variability, particularly over multi-decadal timescales, or due to differences in forced variability compared to the training data. Having completed this method evaluation, our focus now shifts to examining the outcomes when the U-Net is employed with observational data.

549

4.1 Filtering of the observations

The U-Net is now employed to process SAT observations derived from GISSTEMP. 550 By utilizing observed data as input, the U-Net provides an estimate of the forced vari-551 ability. In the interval from 1996 to 2016, the U-Net-derived forced SAT (depicted in Fig-552 ure 8a) illustrates a fairly uniform warming, with amplified warming evident over the 553 Arctic region, consistent with Arctic amplification. Furthermore, this warming effect is 554 slightly more pronounced over land compared to oceans. Conversely, the Southern Ocean 555 experiences less warming in comparison to other global regions. The spatial distribution 556 of standard deviations (Figure 8b), computed from 1905 to 2016 using U-Net output, 557 mirrors the anomalies observed in the 1996-2016 period. This agreement indicates the 558 prevailing influence of increasing anthropogenic forcing. Notably, this pattern closely re-559 sembles the changes observed in the multi-model mean (MMM) (as depicted in Fig. 2). 560 This underscore the significant contribution of the training dataset in determining the 561 identified forced changes. 562

-23-



Figure 7. Spatial average of the RMSE for the forced variability estimated with the U-Net outputs obtained from each ensemble member, and the forced variability obtained with ensemble averages subsampling ensemble of size 1 to 40; for (red) MPI-ESM and (blue) FGOALS-g3. The RMSE calculated from the U-Net and each ensemble member is given by (color shade) the interval including 90% of the distribution, assuming a gaussian distribution, and (horizontal dashed line) the mean RMSE. The RMSE calculated from 500 subsample of size between 1 to 40 is illustrated with (vertical lines) the intervals including 90% of the ensemble member distribution, also assuming a gaussian distribution.



Figure 8. Forced surface air temperature (in °C) anomaly when applying the U-Net to GIS-STEMP observation : a) time average in 1996-2016; b) standard deviation in 1905-2016.



Figure 9. Standard deviation of the SAT deviations from the forced SAT, as estimated using the U-Net, in 1905-2016.

To quantify internal variability within the observations, we compute the deviations 563 of observed SAT anomalies from the estimated forced changes. The resulting internal 564 variability pattern, illustrated by the time standard deviation of these deviations shown 565 in Figure 9, mirrors the model-derived pattern (Fig. 1). Higher internal variability val-566 ues are observed over land areas, as well as regions near the boundaries of sea ice, such 567 as the Labrador Sea and the Nordic Seas in the Northern Hemisphere, and the South-568 ern Ocean. Notably, a local maximum of internal variability emerges in the equatorial 569 Pacific, corresponding to the El Niño-Southern Oscillation region. This similarity in the 570 spatial distribution of internal variability between observations and models underscores 571 the consistency of our findings. 572

We now shift our focus to the GSAT and the Niño 3.4 region (5°N-5°S, 170°W-120°W), 573 with a particular emphasis on Niño 3.4 due to its notably improved performance in our 574 study. In the global context (Figure 10a), the forced variability reveals a consistent warm-575 ing trend, which becomes more pronounced during the 1960s. Notably, the major vol-576 canic eruptions of Agung (1963), El Chichón (1982), and Pinatubo (1991) are associated 577 with temporary cooling patterns. By 2016, the GSAT anomaly reaches 0.7°C. As expected, 578 the forced variability time series exhibits a significant reduction in inter-annual variabil-579 ity. This reduction is particularly striking within the Niño 3.4 region (Figure 10b), where 580 variability at 2 to 7 years is almost entirely eliminated. The U-Net estimates the Niño 581

-26-



Figure 10. Time series of (red) the observed SAT anomaly and (blue) the forced SAT anomaly estimated by the U-Net for a) the global mean b) NINO 3.4 and c) the relative SAT, calculated as the difference between the averaged SAT in Niño 3.4 region and the tropical ocean SAT (30°S-30°N).

3.4 forced variability, depicting a steady warming trend. To quantify the changes of SAT 582 in Niño 3.4 relative to the tropics, we calculate also the relative SAT, defined as the dif-583 ference between the average SAT on the NINO 3.4 region and the average SAT on ocean 584 grid between 30°S-30°N. The relative SST shows that the warming over the Niño 3.4 fol-585 lows that of the tropics, so that no clear El Niño-like reponse is found, unlike climate 586 models (Fig. 2). Some authors (Clement et al., 1996; Heede et al., 2020) have suggested 587 that a forced cooling could exists in the relative SAT, called thermostat effect. Here the 588 relative SAT shows a very small cooling (see Fig. 10c). In addition the SAT in the Niño 589 3.4 region are not affected by the forcing from the main volcanic eruptions. Therefore, 590 no evidence of a Niño-like response to volcanic eruption (as in Khodri et al. (2017)) is 591 found.

592

593 5 Conclusion

A novel approach is introduced in this study to effectively eliminate internal vari-594 ability from a time-evolving two-dimensional dataset, specifically focusing on surface air 595 temperature. The method employs a U-Net neural network and draws inspiration from 596 the noise-to-noise technique. This framework treats internal variability as an analogous 597 noise superimposed on the underlying forced variability. The U-Net model is trained us-598 ing outputs from a diverse ensemble of climate models obtained from the CMIP simu-599 lations. Subsequently, this trained network is applied to observational data to unveil the 600 forced variability signal by attenuating internal variability. The validation of this method 601 involves utilizing large ensemble simulations from individual models, specifically the MPI-602 ESM and FGOALS-g3, to gauge its effectiveness. The forced variability derived from the 603 ensemble mean is then contrasted with the outcomes from the U-Net application. To quan-604 titatively assess the U-Net's efficacy in reducing internal variability, an "equivalent en-605 semble size" is computed. This metric indicates the ensemble size that would be required 606 to achieve the same level of precision in capturing forced changes as the U-Net which is 607 applied to a single member. The U-Net outputs for these two climate models' test data 608 exhibit an error equivalent to an internal variability reduction of a factor of more than 609 4. This magnitude corresponds to the internal variability one could expect from an en-610 semble averaging 17 to 20 members. Furthermore, when the U-Net is applied to surface 611 air temperature observations, the inferred forced changes align closely with the multi-612 model mean in terms of spatial patterns. The U-Net's results do not suggest an El Niño-613 like response to global warming. We observe that the U-Net encounters greater challenges 614 in accurately estimating forced variability over the Arctic region. This discrepancy can 615 be attributed to the significant forced and internal variability associated with changes 616 in sea-ice extent in that area. Additionally, the U-Net's performance in capturing forced 617 variability in the North Atlantic is less successful for the FGOALS-g3 model. This lim-618 itation might be linked to uncertainties stemming from the multi-decadal variability preva-619 lent in these regions (Menary & Wood, 2018; Zhang, 2007). 620

621 622

623

624

625

In the pursuit of enhancing the U-Net methodology, several avenues for future improvements have been identified. One potential approach is to address the U-Net's sensitivity to the multi-model consensus of future variability by employing neural network regularization techniques, such as weights penalisation. Additionally, preprocessing methods like data augmentation could be explored to potentially mitigate such impacts. Im-

-28-

proving the evaluation process of the U-Net's performance is also on the horizon. This 626 could involve testing the U-Net on a broader range of climate models to assess its gen-627 eralizability. Comparing its outcomes with results from alternative methods, such as signal-628 to-noise filtering, could offer a comprehensive evaluation of the U-Net's effectiveness. To 629 broaden the scope of application, the U-Net's performance might be further investigated 630 using additional climate variables beyond surface air temperature (SAT). Variables such 631 as sea level surface pressure and precipitation could be explored, capitalizing on poten-632 tial correlations among these variables to provide more comprehensive insights. Lastly, 633 the proposed method holds the potential for wider applications, including its deployment 634 on simulations from projects like the Detection and Attribution Model Intercomparison 635 Project (Gillett et al., 2016) or the Large Ensemble Single Forcing Model Intercompara-636 ison Project (D. M. Smith et al., 2022). By leveraging transfer learning, the U-Net trained 637 on historical simulations could be adapted to these datasets. This adaptation could fa-638 cilitate the evaluation of specific forcing effects in individual climate models, offering a 639 valuable tool for studying the impact of different external factors on the climate system. 640 Such extensions of the method could contribute significantly to our understanding of cli-641 mate attribution and variability. 642

643 Acknowledgments

We acknowledge the support of the SCAI doctoral program managed by the ANR with
the reference ANR-20-THIA-0003, the support of the EUR IPSL Climate Graduate School
project managed by the ANR under the "Investissements d'avenir" programme with the
reference ANR-11-IDEX-0004-17-EURE-0006. This work was performed using HPC resources from GENCI-TGCC A0090107403 and A0110107403, and GENCI-IDRIS AD011013295.
Guillaume Gastineau was funded by the JPI climate/JPI ocean ROADMAP project (grant
number ANR-19-JPOC-003).

651 6 Open Research

652

Data Availability Statement

The CMIP5 and CMIP6 data is available through the Earth System Grid Feder-

ation and can be accessed through different international nodes. For example : https://

esgf-node.ipsl.upmc.fr/projects/esgf-ipsl/

656 657

677

Codes used in this article for the backward optimization and the figures are from Bône (2023) software available freely at https://zenodo.org/record/8233743.

658 References

- Allen, M. R., & Stott, P. A. (2003). Estimating signal amplitudes in optimal finger printing, Part I: Theory. *Climate Dynamics*, 21, 477–491.
- Allen, M. R., & Tett, S. F. (1999). Checking for model consistency in optimal finger printing. *Climate Dynamics*, 15, 419–434.
- Bonnet, R., Boucher, O., Deshayes, J., Gastineau, G., Hourdin, F., Mignot, J., ...
- Swingedouw, D. (2021). Presentation and evaluation of the ipsl-cm6a-lr ensemble of extended historical simulations. Journal of Advances in Modeling Earth
 Systems, 13(9), e2021MS002565.
- Bonnet, R., Boucher, O., Vrac, M., & Jin, X. (2022). Sensitivity of bias adjustment
 methods to low-frequency internal climate variability over the reference period:
 an ideal model study. *Environmental Research: Climate*, 1(1), 011001.
- Bône, C. (2023). Codes for "Separation of internal and forced variability of climate
 using a U-Net" [Software]. Retrieved from https://zenodo.org/record/
 8233743
- Cassou, C., Kushnir, Y., Hawkins, E., Pirani, A., Kucharski, F., Kang, I.-S., &
 Caltabiano, N. (2018). Decadal Climate Variability and Predictability: Challenges and Opportunities. Bulletin of the American Meteorological Society,
 99(3), 479 490. Retrieved from https://journals.ametsoc.org/view/

journals/bams/99/3/bams-d-16-0286.1.xml

- ⁶⁷⁸ Chylek, P., Li, J., Dubey, M., Wang, M., & Lesins, G. (2011). Observed and model
 ⁶⁷⁹ simulated 20th century Arctic temperature variability: Canadian earth system
 ⁶⁸⁰ model CanESM2. Atmospheric Chemistry and Physics Discussions, 11(8),
 ⁶⁸¹ 22893–22907.
- ⁶⁶² Clement, A. C., Seager, R., Cane, M. A., & Zebiak, S. E. (1996). An ocean dynami ⁶⁶³ cal thermostat. *Journal of Climate*, 9(9), 2190–2196.
- Collier, M. A., Jeffrey, S. J., Rotstayn, L. D., Wong, K., Dravitzki, S., Moseneder,
- C., ... others (2011). The CSIRO-Mk3. 6.0 Atmosphere-Ocean GCM: partici pation in CMIP5 and data publication. In *International congress on modelling and simulation-modsim* (pp. 2691–2697).

688	Collins, W. D., Rasch, P. J., Boville, B. A., Hack, J. J., McCaa, J. R., Williamson,
689	D. L., Zhang, M. (2006). The formulation and atmospheric simulation of
690	the Community Atmosphere Model version 3 (CAM3). Journal of Climate,
691	19(11), 2144-2161.
692	DelSole, T., Tippett, M. K., & Shukla, J. (2011). A significant component of un-
693	forced multidecadal variability in the recent acceleration of global warming.
694	Journal of Climate, 24(3), 909–926.
695	Deser, C., Lehner, F., Rodgers, K. B., Ault, T., Delworth, T. L., DiNezio, P. N.,
696	\dots others (2020). Insights from Earth system model initial-condition large
697	ensembles and future prospects. Nature Climate Change, $10(4)$, 277–286.
698	Deser, C., Phillips, A., Bourdette, V., & Teng, H. (2012). Uncertainty in climate
699	change projections: the role of internal variability. Climate dynamics, $38, 527-$
700	546.
701	Deser, C., & Phillips, A. S. (2023). A range of outcomes: the combined effects of
702	internal variability and anthropogenic forcing on regional climate trends over
703	Europe. Nonlinear Processes in Geophysics, $30(1)$, 63–84.
704	Deser, C., Phillips, A. S., Alexander, M. A., & Smoliak, B. V. (2014). Projecting
705	North American climate over the next 50 years: Uncertainty due to internal
706	variability. Journal of Climate, 27(6), 2271–2296.
707	Döscher, R., Acosta, M., Alessandri, A., Anthoni, P., Arneth, A., Arsouze, T.,
708	others (2021) . The EC-earth3 Earth system model for the climate model in-
709	tercomparison project 6. Geoscientific Model Development Discussions, 2021,
710	1 - 90.
711	Egmont-Petersen, M., de Ridder, D., & Handels, H. (2002). Image processing with
712	neural networks—a review. Pattern recognition, $35(10)$, 2279–2301.
713	Enfield, D. B., & Cid-Serrano, L. (2010). Secular and multidecadal warmings in the
714	North Atlantic and their relationships with major hurricane activity. Interna-
715	tional Journal of Climatology: A Journal of the Royal Meteorological Society,
716	30(2), 174-184.
717	England, M. H., McGregor, S., Spence, P., Meehl, G. A., Timmermann, A., Cai,
718	W., Santoso, A. (2014). Recent intensification of wind-driven circulation
719	in the Pacific and the ongoing warming hiatus. Nature climate change, $4(3)$,
720	222-227.

721	Eyring, V., Bony, S., Meehl, G. A., Senior, C. A., Stevens, B., Stouffer, R. J., &
722	Taylor, K. E. (2016). Overview of the Coupled Model Intercomparison Project
723	Phase 6 (CMIP6) experimental design and organization. Geoscientific Model
724	$Development, \ 9(5), \ 1937-1958.$
725	Eyring, V., Gillett, N., Rao, K. A., Barimalala, R., Parrillo, M. B., Bellouin, N.,
726	Zho, B. (2021). Human Influence on the Climate System. In Climate Change
727	2021: The Physical Science Basis. Contribution of Working Group I to the
728	Sixth Assessment Report of the Intergovernmental Panel on Climate Change.
729	Cambridge University Pres.
730	Frankcombe, L. M., England, M. H., Mann, M. E., & Steinman, B. A. (2015). Sep-
731	arating internal variability from the externally forced climate response. $Journal$
732	of Climate, $28(20)$, $8184-8202$.
733	Frankignoul, C., Gastineau, G., & Kwon, YO. (2017). Estimation of the SST
734	response to anthropogenic and external forcing and its impact on the Atlantic
735	multidecadal oscillation and the Pacific decadal oscillation. Journal of Climate,
736	30(24), 9871-9895.
737	Fyfe, J. C., Kharin, V. V., Santer, B. D., Cole, J. N., & Gillett, N. P. (2021). Sig-
738	nificant impact of forcing uncertainty in a large ensemble of climate model
739	simulations. Proceedings of the National Academy of Sciences, 118(23),
740	e2016549118.
741	Gent, P. R., Danabasoglu, G., Donner, L. J., Holland, M. M., Hunke, E. C., Jayne,
742	S. R., others (2011). The community climate system model version 4.
743	Journal of climate, 24(19), 4973–4991.
744	Gillett, N. P., Shiogama, H., Funke, B., Hegerl, G., Knutti, R., Matthes, K.,
745	Tebaldi, C. (2016). The detection and attribution model intercomparison
746	project (DAMIP v1. 0) contribution to CMIP6. Geoscientific Model Develop-
747	$ment, \ 9(10), \ 3685-3697.$
748	Gulev, S. K., Thorne, P. W., Ahn, J., Dentener, F. J., Domingues, C. M., Gerland,
749	S., others (2021). Changing state of the climate system.
750	Hansen, J., Ruedy, R., Sato, M., & Lo, K. (2010). Global surface temperature
751	change. Reviews of Geophysics, $48(4)$.
752	Harzallah, A., & Sadourny, R. (1995). Internal versus SST-forced atmospheric vari-
753	ability as simulated by an atmospheric general circulation model. Journal of

754	$Climate, \ 8(3), \ 474-495.$
755	Hasselmann, K. (1993). Optimal fingerprints for the detection of time-dependent cli-
756	mate change. Journal of Climate, $6(10)$, 1957–1971.
757	Hawkins, E., & Sutton, R. (2009). The potential to narrow uncertainty in regional
758	climate predictions. Bulletin of the American Meteorological Society, $90(8)$,
759	1095-1108.
760	He, C., Clement, A. C., Cane, M. A., Murphy, L. N., Klavans, J. M., & Fenske,
761	T. M. (2022). A North Atlantic warming hole without ocean circulation.
762	Geophysical research letters, $49(19)$, e2022GL100420.
763	Heede, U. K., Fedorov, A. V., & Burls, N. J. (2020). Time scales and mechanisms
764	for the tropical Pacific response to global warming: A tug of war between the
765	ocean thermostat and weaker Walker. Journal of Climate, $33(14)$, $6101-6118$.
766	Ilesanmi, A. E., & Ilesanmi, T. O. (2021). Methods for image denoising using convo-
767	lutional neural network: a review. Complex & Intelligent Systems, $7(5)$, 2179–
768	2198.
769	Jeffrey, S., Rotstayn, L., Collier, M., Dravitzki, S., Hamalainen, C., Moeseneder, C.,
770	\ldots Syktus, J. (2013). Australia's CMIP5 submission using the CSIRO-Mk3. 6
771	model. Australian Meteorological and Oceanographic Journal, $63(1)$, 1–13.
772	Jiang, W., Gastineau, G., & Codron, F. (2021). Multicentennial variability driven
773	by salinity exchanges between the Atlantic and the Arctic Ocean in a cou-
774	pled climate model. Journal of Advances in Modeling Earth Systems, 13(3),
775	e2020MS002366.
776	Kay, J. E., Deser, C., Phillips, A., Mai, A., Hannay, C., Strand, G., others
777	(2015). The Community Earth System Model (CESM) large ensemble project:
778	A community resource for studying climate change in the presence of internal
779	climate variability. Bulletin of the American Meteorological Society, $96(8)$,
780	1333–1349.
781	Keil, P., Mauritsen, T., Jungclaus, J., Hedemann, C., Olonscheck, D., & Ghosh, R.
782	(2020). Multiple drivers of the North Atlantic warming hole. <i>Nature Climate</i>
783	$Change, \ 10(7), \ 667-671.$
784	Khodri, M., Izumo, T., Vialard, J., Janicot, S., Cassou, C., Lengaigne, M., oth-
785	ers (2017). Tropical explosive volcanic eruptions can trigger El Niño by cooling
786	tropical Africa. Nature communications, $\mathcal{S}(1)$, 778.

787	Kosaka, Y., & Xie, SP. (2013). Recent global-warming hiatus tied to equatorial Pa-
788	cific surface cooling. <i>Nature</i> , 501(7467), 403–407.
789	Lehtinen, J., Munkberg, J., Hasselgren, J., Laine, S., Karras, T., Aittala, M., & Aila,
790	T. (2018). Noise2Noise: Learning image restoration without clean data. $arXiv$
791	preprint arXiv:1803.04189.
792	Lenssen, N. J., Schmidt, G. A., Hansen, J. E., Menne, M. J., Persin, A., Ruedy,
793	R., & Zyss, D. (2019). Improvements in the GISTEMP uncertainty model.
794	Journal of Geophysical Research: Atmospheres, 124(12), 6307–6326.
795	Li, Yu, Y., Tang, Y., Lin, P., Xie, J., Song, M., others (2020). The flexible global
796	ocean-atmosphere-land system model grid-point version 3 (FGOALS-g3): de-
797	scription and evaluation. Journal of Advances in Modeling Earth Systems,
798	12(9), e2019MS002012.
799	Li, S., & Huang, P. (2022). An exponential-interval sampling method for evaluat-
800	ing equilibrium climate sensitivity via reducing internal variability noise. Geo -
801	science Letters, $9(1)$, 1–10.
802	Maher, N., Milinski, S., Suarez-Gutierrez, L., Botzet, M., Dobrynin, M., Kornblueh,
803	L., others (2019) . The Max Planck Institute Grand Ensemble: enabling
804	the exploration of climate system variability. Journal of Advances in Modeling
805	Earth Systems, 11(7), 2050–2069.
806	Marini, C., & Frankignoul, C. (2014). An attempt to deconstruct the Atlantic multi-
807	decadal oscillation. Climate dynamics, 43, 607–625.
808	Masson-Delmotte, V., Zhai, P., Pirani, A., Connors, S. L., Péan, C., Berger, S.,
809	others (2021). Climate change 2021: the physical science basis. Contribution of
810	working group I to the sixth assessment report of the intergovernmental panel
811	on climate change, 2.
812	Meehl, G. A., Hu, A., Arblaster, J. M., Fasullo, J., & Trenberth, K. E. (2013).
813	Externally forced and internally generated decadal climate variability associ-
814	ated with the Interdecadal Pacific Oscillation. $Journal of Climate, 26(18),$
815	7298–7310.
816	Menary, M. B., Robson, J., Allan, R. P., Booth, B. B., Cassou, C., Gastineau, G.,
817	\ldots others (2020). Aerosol-forced AMOC changes in CMIP6 historical simula-
818	tions. Geophysical Research Letters, 47(14), e2020GL088166.
819	Menary, M. B., & Wood, R. A. (2018). An anatomy of the projected north atlantic

-34-
820	warming hole in cmip5 models. Climate Dynamics, $50(7-8)$, $3063-3080$.			
821	Neelin, J. D., Battisti, D. S., Hirst, A. C., Jin, FF., Wakata, Y., Yamagata, T., &			
822	Zebiak, S. E. (1998). ENSO theory. Journal of Geophysical Research: Oceans,			
823	103(C7), 14261-14290.			
824	Newman, M., Alexander, M. A., Ault, T. R., Cobb, K. M., Deser, C., Di Lorenzo,			
825	E., others (2016). The Pacific decadal oscillation, revisited. Journal of			
826	Climate, 29(12), 4399-4427.			
827	O'Shea, K., & Nash, R. (2015). An introduction to convolutional neural networks.			
828	arXiv preprint arXiv:1511.08458.			
829	Parker, D., Folland, C., Scaife, A., Knight, J., Colman, A., Baines, P., & Dong, B.			
830	(2007). Decadal to multidecadal variability and the climate change back-			
831	ground. Journal of Geophysical Research: Atmospheres, 112(D18).			
832	Parsons, L. A., Brennan, M. K., Wills, R. C., & Proistosescu, C. (2020). Magnitudes			
833	and spatial patterns of interdecadal temperature variability in CMIP6. Geo-			
834	$physical\ Research\ Letters,\ 47(7),\ e2019 GL086588.$			
835	Rodgers, K. B., Lin, J., & Frölicher, T. L. (2015). Emergence of multiple ocean			
836	ecosystem drivers in a large ensemble suite with an Earth system model. Bio -			
837	$geosciences, \ 12(11), \ 3301-3320.$			
838	Ronneberger, O., Fischer, P., & Brox, T. (2015). U-net: Convolutional networks for			
839	biomedical image segmentation. In Medical image computing and computer-			
840	$assisted\ intervention-miccai\ 2015:\ 18th\ international\ conference,\ munich,$			
841	germany, october 5-9, 2015, proceedings, part iii 18 (pp. 234–241).			
842	Schmidt, A., Mills, M. J., Ghan, S., Gregory, J. M., Allan, R. P., Andrews, T.,			
843	others (2018). Volcanic radiative forcing from 1979 to 2015. Journal of			
844	Geophysical Research: Atmospheres, 123(22), 12491–12508.			
845	Sherwood, S., Webb, M. J., Annan, J. D., Armour, K. C., Forster, P. M., Har-			
846	greaves, J. C., others (2020). An assessment of Earth's climate sen-			
847	sitivity using multiple lines of evidence. $Reviews of Geophysics, 58(4),$			
848	e2019RG000678.			
849	Smith, C. J., & Forster, P. M. (2021). Suppressed late-20th century warming in			
850	CMIP6 models explained by forcing and feedbacks. Geophysical Research Let-			
851	ters, 48(19), e2021GL094948.			

Smith, C. J., Kramer, R. J., Myhre, G., Alterskjær, K., Collins, W., Sima, A., ...

853	others (2020). Effective radiative forcing and adjustments in CMIP6 models.
854	Atmospheric Chemistry and Physics, 20(16), 9591–9618.
855	Smith, D. M., Gillett, N. P., Simpson, I. R., Athanasiadis, P. J., Baehr, J., Bethke,
856	I., \ldots others (2022). Attribution of multi-annual to decadal changes in the
857	climate system: The Large Ensemble Single Forcing Model Intercomparison
858	Project (LESFMIP). Frontiers in Climate, 4.
859	Solomon, A., Goddard, L., Kumar, A., Carton, J., Deser, C., Fukumori, I., oth-
860	ers (2011). Distinguishing the roles of natural and anthropogenically forced
861	decadal climate variability: implications for prediction. Bulletin of the Ameri-
862	can Meteorological Society, 92(2), 141–156.
863	Steinman, B. A., Mann, M. E., & Miller, S. K. (2015). Atlantic and Pacific mul-
864	tidecadal oscillations and Northern Hemisphere temperatures. Science,
865	347(6225), 988-991.
866	Sun, L., Alexander, M., & Deser, C. (2018). Evolution of the global coupled climate
867	response to Arctic sea ice loss during 1990–2090 and its contribution to climate
868	change. Journal of Climate, 31(19), 7823–7843.
869	Swart, N. C., Fyfe, J. C., Hawkins, E., Kay, J. E., & Jahn, A. (2015). Influence of
870	internal variability on Arctic sea-ice trends. Nature Climate Change, $5(2)$, 86–
871	89.
872	Swingedouw, D., Bily, A., Esquerdo, C., Borchert, L. F., Sgubin, G., Mignot, J., &
873	Menary, M. (2021). On the risk of abrupt changes in the north atlantic sub-
874	polar gyre in cmip6 models. Annals of the New York Academy of Sciences,
875	1504(1), 187-201.
876	Taylor, K. E., Stouffer, R. J., & Meehl, G. A. (2012). An overview of CMIP5 and
877	the experiment design. Bulletin of the American meteorological Society, $93(4)$,
878	485–498.
879	Tebaldi, C., Debeire, K., Eyring, V., Fischer, E., Fyfe, J., Friedlingstein, P., oth-
880	ers (2020). Climate model projections from the scenario model intercomparison
881	project (ScenarioMIP) of CMIP6. Earth System Dynamics Discussions, 2020,
882	1-50.
883	Tian, C., Fei, L., Zheng, W., Xu, Y., Zuo, W., & Lin, CW. (2020). Deep learning
884	on image denoising: An overview. Neural Networks, 131, 251–275.
885	Ting, M., Kushnir, Y., Seager, R., & Li, C. (2009). Forced and internal twentieth-

-36-

886	century SST trends in the North Atlantic. J. Climate, 22, 1469–1481.
887	Van Vuuren, D. P., Edmonds, J., Kainuma, M., Riahi, K., Thomson, A., Hibbard,
888	K., others (2011). The representative concentration pathways: an overview.
889	Climatic change, 109, 5–31.
890	Vincent, L., Zhang, X., Brown, R., Feng, Y., Mekis, E., Milewska, E., Wang, X.
891	(2015). Observed trends in Canada's climate and influence of low-frequency
892	variability modes. Journal of Climate, 28(11), 4545–4560.
893	Wang, C., & Picaut, J. (2004). Understanding ENSO physics—A review. Earth's
894	Climate: The Ocean–Atmosphere Interaction, Geophys. Monogr, 147, 21–48.
895	Wills, R. C., Battisti, D. S., Armour, K. C., Schneider, T., & Deser, C. (2020). Pat-
896	tern recognition methods to separate forced responses from internal variability
897	in climate model ensembles and observations. Journal of Climate, $33(20)$,
898	8693–8719.
899	Yamashita, R., Nishio, M., Do, R. K. G., & Togashi, K. (2018). Convolutional neu-
900	ral networks: an overview and application in radiology. Insights into imaging,
901	9, 611-629.
902	Zelinka, M. D., Myers, T. A., McCoy, D. T., Po-Chedley, S., Caldwell, P. M., Ceppi,
903	P., Taylor, K. E. (2020). Causes of higher climate sensitivity in CMIP6
904	models. Geophysical Research Letters, $47(1)$, e2019GL085782.
905	Zelinka, M. D., Zhou, C., & Klein, S. A. (2016). Insights from a refined decomposi-
906	tion of cloud feedbacks. Geophysical Research Letters, $43(17)$, $9259-9269$.
907	Zhang, R. (2007). Anticorrelated multidecadal variations between surface and sub-
908	surface tropical north atlantic. Geophysical Research Letters, $34(12)$.
909	Zhang, R., Sutton, R., Danabasoglu, G., Kwon, YO., Marsh, R., Yeager, S. G.,
910	Little, C. M. (2019). A review of the role of the Atlantic meridional over-
911	turning circulation in Atlantic multidecadal variability and associated climate
912	impacts. Reviews of Geophysics, 57(2), 316–375.

12

Supporting Information for "Separation of internal and forced variability of climate using a U-Net"

Constantin Bône¹², Guillaume Gastineau¹, Sylvie Thiria¹, Patrick

Gallinari²³and Carlos Mejia¹

¹UMR LOCEAN, IPSL, Sorbonne Université, IRD, CNRS, MNHN

 $^2 \mathrm{UMR}$ ISIR, Sorbonne Université, CNRS, INSERM

 $^3\mathrm{Criteo}$ AI Lab

Contents of this file

- 1. Figures S1 to S2
- 2. Tables S1 to S3

X - 2

 Table 1.
 List of climate CMIP5 climate model used. Nb indicates the ensemble size of each simulation.

:

Table 1.		
Name	Nb Historical	Nb RCP8.5
GISS-E2-R-p1	2	2
MPI-ESM-LR	3	3
CanESM2	5	5
CESM1-CAM5	3	3
FIO-ESM	3	3
CNRM-CM5	10	5
CSIRO-Mk3-6-0	10	10
FGOALS-g3-s2	3	3
GISS-E2-H-p1	6	3
GISS-E2-H-p3	2	2
HadGEM2-ES	3	3
IPSL-CM5A-LR	6	4
GISS-E2-R-p3	3	2

Nb Historical Nb SSP2-4.5 Name ACCESS-CM2 CanESM5 CESM2 CanESM5-CanOE GISS-E2-1-G EC-Earth3 MIROC-ES2L HadGEM3-GC31-LL GFDL-ESM4 FIO-ESM-2-0 KACE-1-0-G GISS-E2-1-G-p3 ACCESS-ESM1-5 CAS-ESM2-0 NESM3 MPI-ESM1-2-HR NorESM2-LM GISS-E2-1-G-p5 IPSL-CM6A-LR GISS-E2-1-H CESM2-WACCM CNRM-CM6-1 CAMS-CSM1-0 UKESM1-0-LL MPI-ESM1-2-LR MRI-ESM2-0 CNRM-ESM2-1 FGOALS-f3 CanESM5 MIROC6

 Table 2.
 List of climate CMIP6 climate model used. Nb indicates the ensemble size of each simulation.

:

X - 4

 Table 3.
 List of climate SMILE climate model used. Nb indicates the ensemble size of each simulation.

Name of the model	Nb of historical members	Nb scenario members	Origin of forcings	Scenario
CSIRO-Mk3-6-0	29	29	CMIP5	RCP8.5
EC-EARTH	15	15	CMIP6	SSP2-4.5
MPI-ESM	100	100	CMIP6	SSP2-4.5
FGOALS-g3	110	110	CMIP6	SSP2-4.5

:



Figure S1. Validation RMSE (in °C) using the ensemble mean of MIROC6 outputs as a target, and each member as inputs for different epochs and when varying the numbers of filters for each convolutionnal layer of the U-Net. Vertical line of the same colour shows the epoch where the minimum RMSE is obtained for the three changes in the number of filters.



Figure S2. a) Spatial average of the RMSE, in °C, between the U-Net output obtained from each member of MPI-ESM and the ensemble average (blue) over 90°S-90°N and (red) over 60°N-90°N. The line provides the ensemble mean error obtained with an average from the errors of all U-Net outputs. Colour shade shows one standard deviation among the error of the outputs from all members. b) Same as a) but for FGOALS-g3 members.