Understanding model-observation discrepancies in satellite retrievals of atmospheric temperature using GISS ModelE

Madeline Claire Casas¹, Gavin A. Schmidt², Ron L. Miller², Clara Orbe², Larissa S. Nazarenko³, and Susanne E. Bauer⁴

¹Stanford School of Humanities and Sciences
²NASA Goddard Institute for Space Studies
³Columbia University/NASA GISS
⁴NASA Goddard Institute for Space Studies, New York, NY, USA

November 23, 2022

Abstract

We examine multiple factors in the representation of satellite-retrieved atmospheric temperature diagnostics in historical simulations of climate change during the satellite era (specifically 1979-2021) using GISS ModelE contributions to the Coupled Model Intercomparison Project (Phase 6) (CMIP6). The tropospheric and stratospheric trends in these diagnostics are affected by greenhouse gases (notably carbon dioxide and ozone), coupling with the ocean, volcanic aerosols, solar activity and compositional and dynamic feedbacks. We explore the impacts of internal variability, changing forcing specifications, composition interactivity, the quality of the stratospheric circulation, vertical resolution, and possible impacts of the mis-specification of volcanic aerosol optical depths.

Overall trends and patterns over the satellite period are well captured, but discrepancies at all levels exist and have multiple distinct causes. We find that stratospheric comparisons (using Stratospheric Sounding Unit (SSU) retrievals and successor instruments) are most affected by variations in the representation of ozone depletion and feedbacks, followed by the volcanic signals. Tropospheric skill (using the Microwave Sounding Unit (MSU) retrievals) is affected by the trends in ocean temperature and tropospheric aerosols, but also by the representation of stratospheric processes through the impact of the Brewer-Dobson circulation on the height of the tropical tropopause. We do not find evidence of a systematic problem in the model climate sensitivity.

Understanding model-observation discrepancies in satellite retrievals of atmospheric temperature using GISS ModelE

Madeline C. Casas^{1,2}, Gavin A. Schmidt¹, Ron. L. Miller¹, Clara Orbe¹, Larissa S. Nazarenko^{1,3}, and Susanne E. Bauer¹

 $^1\mathrm{NASA}$ Goddard Institute for Space Studies, New York, N
Y $^2\mathrm{Stanford}$ School of Humanities and Sciences, Palo Alto, C
A $^3\mathrm{Center}$ for Climate System Research, Columbia University, New York, NY

Key Points:

10	•	Changes in forcing configurations have significant impact on agreement with satel-
11		lite data.
12	•	Tropospheric model-observation agreement is linked to ocean heat uptake, ozone
13		and internal variability.
14	•	Stratospheric discrepancies are related to volcanic aerosol modeling and ozone sen-

sitivity.

1

2

3

4

5

6 7 8

9

15

Corresponding author: Gavin A. Schmidt, gavin.a.schmidt@nasa.gov

16 Abstract

We examine multiple factors in the representation of satellite-retrieved atmospheric tem-17 perature diagnostics in historical simulations of climate change during the satellite era 18 (specifically 1979–2021) using GISS ModelE contributions to the Coupled Model Inter-19 comparison Project (Phase 6) (CMIP6). The tropospheric and stratospheric trends in 20 these diagnostics are affected by greenhouse gases (notably carbon dioxide and ozone), 21 coupling with the ocean, volcanic aerosols, solar activity and compositional and dynamic 22 feedbacks. We explore the impacts of internal variability, changing forcing specifications, 23 composition interactivity, the quality of the stratospheric circulation, vertical resolution, 24 and possible impacts of the mis-specification of volcanic aerosol optical depths. 25

Overall trends and patterns over the satellite period are well captured, but discrep-26 ancies at all levels exist and have multiple distinct causes. We find that stratospheric 27 comparisons (using Stratospheric Sounding Unit (SSU) retrievals and successor instru-28 ments) are most affected by variations in the representation of ozone depletion and feed-29 backs, followed by the volcanic signals. Tropospheric skill (using the Microwave Sound-30 ing Unit (MSU) retrievals) is affected by the trends in ocean temperature and tropospheric 31 aerosols, but also by the representation of stratospheric processes through the impact 32 of the Brewer-Dobson circulation on the height of the tropical tropopause. We do not 33 find evidence of a systematic problem in the model climate sensitivity. 34

³⁵ Plain Language Summary

The assessment of the ability of climate models to match observed trends and vari-36 ability seen in the real world is a key factor in building the credibility of their projec-37 tions under future scenarios. We focus on the trends in the atmosphere temperatures from 38 the surface to the stratosphere whose trends at different levels reflect different processes 39 and drivers. The satellite retrievals are weighted averages of atmospheric temperatures 40 and so the vertical structure of the model trends matter in the comparison. We find that 41 overall the trends throughout the atmosphere are well-captured by the GISS models but 42 that discrepancies can occur due to misspecified forcings, internal variability, and model 43 structure. 44

45 **1** Introduction

The start of the satellite period (nominally 1979) marked the dawn of a new era 46 in global multi-variate monitoring of climate. Over 40 years of data have been collected 47 since then and has been sufficient not only to refine our knowledge of the Earth's clima-48 tology, but also to capture the trends of a changing climate. One suite of important vari-49 ables are the vertically-weighted atmospheric temperature changes seen by the Microwave 50 Sounding Units (MSUs) (Spencer & Christy, 1990; Mears et al., 2003), Stratospheric Sound-51 ing Units (SSUs) (Thompson et al., 2012) and their successors, the Advanced Microwave 52 Sounding Units (AMSUs). The clear trends at the surface and in both the troposphere 53 and stratosphere have long been used in detailed comparisons to climate model simu-54 lations. Those comparisons have led to the discovery of discrepancies between the satel-55 lite retrievals, surface temperatures and models, and therefore understanding them has 56 been a focus of scientific attention for more than two decades (e.g. Christy & Spencer, 57 1995; Jones et al., 1997; Hansen et al., 1995; CCSP, 2006; Thorne et al., 2011). 58

The reasons for any climate model's mismatch to these trends can arise from multiple factors: errors in the model physics; model drivers; observational retrieval errors; or simply from an inappropriate comparison, and all of these possible effects have been encountered over time with respect to the MSU/SSU/AMSU time-series (for instance, Wentz & Schabel, 1998; Santer et al., 1999; Fu et al., 2004; Thompson et al., 2012; Santer et al., 2014; Zou & Qian, 2016). With respect to issues related to the retrievals themselves, there have been multiple updates to the various independent products that have
progressively dealt with issues in calibrations, orbits, overlaps, diurnal cycle adjustments
etc. (for instance Mears et al., 2003, 2012; Spencer et al., 2017; Zou & Qian, 2016). Comparisons to the multi-model ensembles have been updated as a consequence (Maycock
et al., 2018; Seidel et al., 2016), reducing some of the differences, but not all, and not
with all observational products.

Mitchell et al. (2020) recently compared model trends at specific heights and found 71 little to no improvement in the model ensemble skill as a whole in going from Phase 5 72 73 of the Coupled Model Intercomparison Project (CMIP5) (circa 2011) to the 6th phase (CMIP6) (2019–2021). Trends in the tropical mid-troposphere still seem too large com-74 pared to radiosonde trends (McKitrick & Christy, 2020). However, both of these papers 75 only looked at a single ensemble member from each model and so their results may be 76 biased by not looking at the full range of internal variability (Po-Chedley et al., 2021). 77 Within those ensembles, however, there are a number of more structured tests that can 78 help illuminate the reasons for the continued discrepancies. In particular, controlled vari-79 ations of model structure, initial conditions, forcings and components within a single model 80 family can be used to examine reasons for the remaining discrepancies. 81

In this paper, we look at the GISS ModelE2.x family of contributions to CMIP6 82 (Table 1). Model configurations include variations in the ocean component (either ob-83 served sea surface temperatures or two different ocean models), the radiative forcings 84 applied, the interactivity of atmospheric composition, vertical resolution and the qual-85 ity of the stratospheric representation (Table 2). Each configuration has multiple ensem-86 ble members (either 5 or 10 members). Additionally, we make use of a suite of single forc-87 ing ensemble experiments (5 members each) that highlight the vertically varying finger-88 prints of specific forcings to help diagnose the changes. 89

1.1 Background

90

The vertical pattern of temperature change in response to increased greenhouse gases 91 has been recognised as a distinct fingerprint since the pioneering work of Manabe and 92 Wetherald (1967). The surface warming, enhanced tropospheric warming, and strato-93 spheric (and above) cooling is unlike the pattern generated by increasing solar activity 94 (which would have more uniform warming through the whole atmosphere), ozone deple-95 tion, volcanic activity, or ocean-driven internal variability. However, it wasn't until the 96 satellite era and the development of global atmospheric retrievals using the MSU/SSU/AMSU 97 series of instruments combined with the radiosonde record that the ability to distinguish 98 these vertical fingerprints emerged (Spencer & Christy, 1990; Randel & Cobb, 1994; San-99 ter et al., 1996; Ramaswamy et al., 1996). 100

The use of these datasets for the detection and attribution of climate change has 101 been complicated by the substantial structural uncertainty associated with the retrievals 102 themselves (Hansen et al., 1995; Mears et al., 2003; CCSP, 2006; Thompson et al., 2012; 103 Po-Chedley & Fu, 2012; Zou & Qian, 2016), though as the trend signal has grown and 104 as successive non-climatic influences have been dealt with, those differences have become 105 less relevant. Nonetheless, the differences in atmospheric trends between models and ob-106 servations continue to generate substantial discussion (McKitrick & Christy, 2020; Mitchell 107 et al., 2020; Fyfe et al., 2021). 108

In the papers referenced above, the structural uncertainty in models is often assessed through a sampling of the CMIP ensembles (over many generations of this project). This is a good way to assess some aspects of that variance - for instance, with respect to the treatment of convection, or the sensitivity to variations in climate sensitivity, but the use of an 'ensemble of opportunity' is not a complete assessment of uncertainty, and some real aspects of the uncertainty will not be sampled at all. Within a single model or model family however, we can address some structural variations in a more controlled

way and specifically address different sources of uncertainty. Specifically, how sensitive 116 are the comparisons to the real uncertainties in the forcing functions? Or to included 117 interactive composition? We have deliberately added these kinds of model variations to 118 the CMIP6 archive, but note that many papers examining the CMIP6 multi-model en-119 semble will only use a single model from a particular model family and sometimes only 120 a single ensemble member. This leads to the potential conflation of internal variability 121 with structural variability, underestimating both, and doesn't take advantage of more 122 controlled variations with model families. Thus, single model family analyses should be 123 seen as complementary and orthogonal to analyses that use the multi-model ensemble. 124

2 Observational Data Sources

Multiple groups have independently analysed the raw MSU, SSU and AMSU data 126 retrievals to create time-series of atmospheric temperatures, notably the University of 127 Alabama Huntsville (UAH) group and Remote Sensing Systems (RSS) (for the MSU data), 128 and the NOAA Center for SaTellite Applications and Research (NOAA STAR) (both 129 MSU and SSU products). We use the latest versions that are publicly accessible (UAH 130 v6, RSS v4, NOAA STAR v4.1 (MSU) and v3.0 (SSU)). We use the differences between 131 them as an indication of the structural uncertainty in the retrieved trends, though we 132 recognise this is possibly an underestimate. The structural uncertainty of the SSU di-133 agnostics is unclear, though reduced compared to previous versions (Thompson et al., 134 2012). We also focus on the global means, with the understanding that the vertical sig-135 nal of trend variability in the troposphere is dominated by the moist-lapse-rate controlled 136 tropical regions. All datasets are used through to the end of 2021, except for the NOAA 137 STAR SSU products for which data only through to the end of 2020 is currently avail-138 able (as of April 2022). 139

In comparing Surface Air Temperature (SAT) trends in the models to the obser-140 vations, we are mindful that trends in a blended product of SST and SAT anomalies (the 141 Land-Ocean Temperature Index (LOTI)) such as produced by GISTEMP, HadCRUT5 142 or the (pending) NOAAGlobalTemp Interim product (Lenssen et al., 2019; Vose et al., 143 2021; Morice et al., 2021) may be systematically different from the pure SAT trends (Richardson 144 et al., 2018). For instance, in the 10 simulations with the GISS E2.1-G f2 configuration 145 (see below), the global mean SAT trends are 0.036 [0.028,0.044] °C/dec (95% range) greater 146 than the LOTI trends for the time period 1979–2014. Thus as an alternative measure, 147 we also use the SAT trends from the European Centre for Medium Range Weather Fore-148 casts (ECMWF) Reanalysis version 5 (ERA5) (Hersbach et al., 2020; Simmons et al., 149 2021), which is perhaps a more appropriate comparison, although in practice it is sim-150 ilar. 151

¹⁵² **3 GISS ModelE simulations**

We analyse model simulations performed for the Coupled Model Intercomparison 153 Project (Phase 6) (CMIP6) using various configurations of GISS ModelE, namely GISS-154 E2.1-G, GISS-E2.1-H (Kelley et al., 2020) and GISS-E2.2-G (Rind et al., 2020). The E2.1 155 model is an update of the GISS-E2 simulations that were used in CMIP5 (Schmidt et 156 al., 2014; Miller et al., 2014) with the same basic resolution $(2.5^{\circ} \times 2^{\circ})$ in the atmosphere, 157 $\approx 1^{\circ}$ in the ocean), but with multiple fixes and improvements in tuning. The -G and 158 -H versions differ in the ocean model while the AMIP simulations use SST as a bound-159 ary forcing (PCMDI-AMIP-v1.1, based on HadISSTv1.1) (Taylor et al., 2000). The E2.2 160 versions have a higher model top (0.002 hPa compared to 0.1 hPa) and twice the ver-161 tical resolution of E2.1 in the atmosphere and have been designed to greatly improve strato-162 spheric circulation and variability. There are also some atmospheric retunings that were 163 made that affect the base climatology and variability (Orbe et al., 2020), notably there 164



Main drivers over the satellite era

Figure 1. Schematic of the important drivers of atmospheric change over the satellite period. Each variable is plotted as a normalized index (with zero mean and unit standard deviation over the period 1979–2021) in order to highlight the pattern of variance over time. Data sources: so-lar irradiance (NRLTSI2) (Coddington et al., 2016), ozone hole area (Kramarova et al., 2014), volcanic stratospheric aerosol optical depth (Sato et al., 1993), well-mixed greenhouse gases and tropospheric aerosol radiative forcing (from the E2.1-G f2 simulations) (Miller et al., 2021). The vertical dotted line distinguishes the 'ozone depletion' and 'ozone recovery' periods for the stratospheric analyses.

is an overall cold bias but a more realistic spectrum and magnitude of ENSO variabil ity.

Each model configuration has options for the interactivity of atmospheric compo-167 sition (specifically gas phase chemistry and aerosol physics). The versions denoted physics_version=1 168 (p1), (NINT) have non-interactive composition, with three-dimensional seasonality (monthly 169 means) and trends in radiatively active components (ozone and aerosols) taken from the 170 interactive physics_version=3 (p3) versions that use the One-Moment Aerosol (OMA) 171 scheme and whole-atmosphere chemistry (Bauer et al., 2020). The aerosol number con-172 173 centrations that impact clouds are obtained from the aerosol mass (Menon & Rotstayn, 2006). Additionally, physics_version=5 (p5) uses the MATRIX aerosol module with 174 the same chemistry (Bauer et al., 2008). In MATRIX the number of cloud activating par-175 ticles is based on an aerosol activation parameterization which treats multimodal and 176 multicomponent aerosols and provide the activated fraction of the number and mass con-177 centration for each population, based on the population composition and the cloud up-178 draft velocity (Abdul-Razzak et al., 1998; Abdul-Razzak & Ghan, 2000). GISS-E2.1 in-179 cludes only the first indirect effect, which is the effect of aerosols on cloud droplet num-180 ber concentration and thereby on cloud albedo, cloud effective radii and radiation (Menon 181 et al., 2008, 2010). Miller et al. (2021) has a fuller description of the differences among 182 the physics versions. 183

We focus on the 'historical' simulations (from 1850–2014) driven by a suite of cli-184 mate drivers, the Shared Socio-Economic Pathway (SSP) scenario 2-4.5 (ssp245) runs 185 (from 2015 onward) (Nazarenko et al., 2022), and supplemented by various single-forcing 186 simulations for the historical period (Fig. 1). We have more simulations for the histor-187 ical period than for the SSPs, but because of the vagaries of El Niño/La Niña cycles, only 188 looking at the trends to 2014 might bias the comparisons. Thus where we have config-189 urations that were run for SSP2-4.5, we also track trends over the longer period. The 190 varying composition forcings for E2.1 used in the non-interactive (NINT) cases are de-191 rived from our interactive (OMA) runs. The NINT f1 forcings came from our initial AMIP 192 runs with E2.1 (OMA). However, the discovery of an error in the coding for stratospheric 193 ozone chemistry led us to later rerun these simulations to generate the f2 suite of forc-194 ings (which differ mainly in the ozone trends in the stratosphere) (Miller et al., 2021). 195 These f2 runs also have a complete suite of results with individual forcings to 2014, with 196 some simulations going to 2018 or continued using the SSP2-4.5 drivers. Finally, we have 197 a set of forcings f3 that are interpolated from the higher vertical resolution E2.2-G (OMA) 198 model that had a noticeably improved stratospheric-tropospheric exchange and circu-199 lation, which impacted the ozone climatology, variability and trends. Note that changes 200 in stratospheric water vapor associated with solar-cycle related photolysis changes are 201 not included in any of the NINT runs. In total, we examine nine separate configurations, 202 with over 50 individual simulations. 203

The ozone and aerosol forcings in each individual configuration may thus be different from the schematic in Figure 1, but the overall transient pattern is close. The MA-TRIX runs have a faster decline of the magnitude of the tropospheric aerosol effects than those in which the aerosols are derived from the OMA runs (Bauer et al., 2020). Similarly, the exact timing of the ozone hole and stabilization is different in the E2.2 models than in E2.1 (Orbe et al., 2020). We will return to these issues in the discussion.

There is a subtle difference between the net anthropogenic forcing in the E2.1-G f3 (NINT) runs and the E2.2-G (OMA) runs, from which the ozone and aerosols were derived, related to the indirect aerosol effect. In the non-interactive composition model configurations, the aerosol indirect effects are tuned so that year 2000 forcing, given the aerosol distribution, is around -1 W/m2 (Miller et al., 2021). This tuning is slightly different for the E2.2 and E2.1 model configurations. Thus, when taking the aerosol distribution from E2.2 and using it in the E2.1 model, there is a difference in the aerosol

Model version	Experiment	ripf number	DOI
E2.1-G	amip	r[1-5]i1p1f2	10.22033/ESGF/CMIP6.6984
	historical	r[1-10]i1p[135]f[123]	10.22033/ESGF/CMIP6.7127
	ssp245	r[1-10]i1p[135]f2	10.22033/ESGF/CMIP6.7415
	hist-volc	r[1-5]i1p1f2	10.22033/ESGF/CMIP6.7111
	hist-sol	r[1-5]i1p1f2	10.22033/ESGF/CMIP6.7101
	hist-aer	r[1-5]i1p1f2	10.22033/ESGF/CMIP6.7081
	hist-GHG	r[1-5]i1p1f2	10.22033/ESGF/CMIP6.7079
	hist-totalO3	r[1-5]i1p1f2	N/A
E2.1-H	historical	r[1-5]i1p1f2	10.22033/ESGF/CMIP6.7128
	ssp245	r[1-5]i1p1f2	10.22033/ESGF/CMIP6.7416
E2.2-G	historical	r[1-5]i1p[13]f1	10.22033/ESGF/CMIP6.6951
	ssp245	r[1-5]i1p3f1	10.22033/ESGF/CMIP6.7415

Table 1. Model experiments in CMIP6, simulation identifiers (using standard regular expression format) and DOIs for the ensemble. The **p** variable denotes different treatment of atmospheric composition, with **p1** being non-interactive (NINT), **p3** using whole atmospheric chemistry and the One Moment Aerosol (OMA) module, and **p5** which uses whole atmosphere chemistry and the MATRIX aerosol scheme (Bauer et al., 2020). Note that the definition of the forcing variants are unique to each model and physics variant (so **f1** in the **p1f1** (NINT) simulations is not related to the **f1** in the **p5f1** (MATRIX) simulations). The hist-totalO3 simulations were not submitted as part of the CMIP6 request, but are included in this analysis for completeness.

indirect effect that leads to a decrease in net forcing of 0.27 W/m^2 compared to the f2 configuration.

We analyse the surface air temperatures (SAT), the Temperature of the Lower Tro-219 posphere (TLT) (3km/700 hPa) the Temperature of Mid-Troposphere (TMT)(5 km/500 220 hPa), the Temperature of the Lower Stratosphere (TLS) (18 km/80 hPa), and SSU chan-221 nels 1, 2, and 3 (centered on 31, 39, 45 km and 10, 3 and 1.5 hPa respectively). Height 222 and pressures are given for the peak in the atmospheric weighting, but the tails of the 223 weighting functions are quite broad and extend over a wide vertical range, necessitat-224 ing an appropriate weighted diagnostic in the models for comparison. The MSU and SSU 225 diagnostics within the model are based on a fixed weighting in pressure and, although 226 more complicated forward models can be applied (Shah & Rind, 1995), they do not no-227 ticeably impact the global trends (Schmidt et al., 2006). 228

Santer et al. (2021) (following Fu et al. (2011)) analysed a version of TMT that uses the TLS to correct for differing estimates of lower stratospheric cooling. They found that the trends in the corrected TMT in CMIP5 and CMIP6 models were not statistically distinct from the TLT trends, and so we choose not to additionally analyse the corrected TMT product.

234 4 Methods

We focus on the annual global anomalies and linear trends in the ensembles for each of the diagnostics described above over the 1979–2014 or 1979–2020/2021 periods. Where needed, we reference anomalies to a 1980–1999 baseline. Ensemble spread is denoted using a 95% confidence interval derived from the 5 or 10 ensemble members.

Using an ordinary least squares linear regression on the annual anomaly data, we calculated the trends in °C per decade for each run and for each variable in the ensem-

Model Configuration	SAT ($^{\circ}C$)	$ECS (^{o}C)$	$TCR (^{o}C)$	Model Top/Layers
E2.1-G f1 (NINT)	14.3	2.7	1.8	0.1 hPa/40 L
E2.1-G f2 (NINT)	14.1	2.7	1.8	н
E2.1-G f3 (NINT)	14.2	2.7	1.8	н
E2.1-G (OMA)	14.7	2.6	1.6	н
E2.1-G (MATRIX)	14.8	2.8	1.8	н
E2.1-H f2 (NINT)	14.5	3.1	1.9	н
E2.2-G (NINT)	12.3	2.4	1.7	0.002 hPa/102 L
E2.2-G (OMA)	11.7	2.1	1.4	ш
Observations	$14.3 {\pm} 0.5$	2.0 - 5.0	1.2 - 2.4	

Table 2. Selected model characteristics for the different configurations used. Global Mean Surface Air Temperature (SAT, °C) is for the period 1981–2010. ECS and TCR are calculated from the abrupt4xCO2 and 1pctCO2 experiments, respectively. Observations are inferred from Jones et al. (1999), and the 'very likely' sensitivity ranges from the IPCC AR6 report (Masson-Delmotte et al., 2021).

ble and for the ensemble mean. Where relevant, uncertainties are given as the 95% confidence interval on the linear regression on the annual data. We construct density plots using the computed decadal trends for each run in an ensemble using the density function in R, following the methods of Sheather and Jones (1991). These plots are used to visualize the spread of decadal trend values within the ensemble and to compare various ensembles for comparison to the observational products for each metric.

In an effort to isolate the effects of stratospheric ozone depletion, we separate our calculation of trends in the stratosphere between the ozone depletion era (1979–1998) and the recovery period (after 1999) following Mitchell et al. (2013) and Seidel et al. (2016). For the TLS data in particular, a single linear trend is not a good fit, and so the separation of these periods can be used to more clearly distinguish the impact of the ozone depletion, as can the inclusion of volcanic and solar predictors in a multiple linear regression.

Consistency of the trends with observational products is assessed in two ways. First, we perform a simple test of whether the ensemble mean from the model configuration is within the 95% confidence interval from one of the observational product(s). This tests whether the observed trend is consistent with our estimate of the forced signal. A better test is whether an observational trend could be plausibly drawn from the model distribution of forced signal plus internal variability. This statistic is calculated following Eqn. 12 in Santer et al. (2008) (assuming a single model) using

$$d = |\overline{T_m} - \overline{T_o}| / \sqrt{s\{\langle T_m \rangle\}^2 + s\{T_o\}^2} \tag{1}$$

where $s\{\langle T_m \rangle\}$ is the standard deviation of the ensemble of model temperature trends 261 T_m and $s\{T_o\}$ is the standard deviation in the linear regression the observed tempera-262 tures, respectively. This d statistic is assumed to follow a Student's t-distribution, and 263 so the probability of getting as high a value as d can be assessed (assuming the degrees 264 of freedom are one less than the ensemble size). If the probability is less than 95% in a 265 two-tailed test, we conclude that the observations are consistent with the specific model 266 ensemble. Since we are using annual data to compute the trends, the degree of autocor-267 relation in the residuals is small and neglected here. Inclusion of this effect would lead 268

to slightly broader confidence intervals, and slightly greater consistency, but this does not impact the pattern of results we see nor any conclusions.

In the troposphere, models and observations indicate that the ratio of tropospheric trends to the SAT trends (related to the effective lapse rate) is more stable than the trends themselves (Wigley, 2006; Santer et al., 2005, 2017). Thus we also examine the ratios of TLT and TMT data to the SAT products in each configuration.

By contrasting specific pairs or sets of simulations in our archive, we can isolate 275 many different aspects of the drivers and responses. For instance, in the troposphere, 276 we can distinguish the impacts of changes in sea surface temperature (SST) (via differ-277 ent ocean models or observed ocean temperatures) by comparing the E2.1-G f2, E2.1-278 H f2 and E2.1 (AMIP) simulations. One such difference arises by varying rates and struc-279 ture of ocean heat uptake in E2.1-H f2 compared to E2.1-G f2. There is more overall 280 ocean heat uptake in E2.1-H f2, however the uptake is localized to the upper ocean rather 281 than the deep ocean. This creates a larger SST increase in E2.1-H f2 than in E2.1-G f2 282 simulations with identical forcings (but see Miller et al. (2021) for a more thorough ex-283 ploration of the ocean heat content changes in the GISS E2.1 simulations). We can also 284 compare E2.2-G to E2.1-G f3 and f2 to examine whether there is significant improve-285 ment related to the higher vertical resolution model and better resolved stratosphere, 286 and whether changes between simulations relates to the forcings or model structure. We 287 have multiple realizations of the forcing fields (notably, aerosols and ozone) with the same 288 underlying climate model to examine the sensitivity to those fields. Also, within each 289 ensemble, we can estimate the impact of the modeled internal variability on the trends. 290

The impact of specific large volcanic eruptions in the first half of the satellite period (El Chichon in 1982 and Mt. Pinatubo in 1991) could bias the trend comparisons if there are issues in either the volcanic forcings used in the models or the model radiative response to the volcanic aerosol inputs. Similarly, there is clear evidence of a solar cycle signal in the stratospheric diagnostics. We therefore use stratospheric volcanic aerosol depth and solar irradiance estimates in an additional multiple linear regression to reduce the influence of potential errors in natural forcings and/or response.

298 5 Results

299

5.1 Tropospheric Trends

We first examine the time-series from a subset of the configurations in figure 2. This 300 shows not only the contrasting trends, but also the degree of simulated internal variabil-301 ity. As expected, the AMIP configurations (driven with observed SST) have very little 302 spread and are a very good match to the RSS and NOAA STAR TLT and TMT changes, 303 though they diverge from UAH TLT, notably after 2000. However, we should note that 304 the AMIP SST files are based on HadSST2 (Rayner et al., 2006). Recent revisions (to 305 HadSST4, (Kennedy et al., 2019)) have dealt with many non-climate discontinuities, and 306 the net effect has been to increase the reported trends - by about 50% in the tropics and 307 globally over the 1979–2019 period. These changes will have an impact on future AMIP-308 style runs (Flannaghan et al., 2014), likely increasing the tropospheric SAT and TLT trends. 309 Assessing the importance of these changes will be the subject of future study. 310

The various flavors of E2.1-G coupled models have more spread due to over-estimated magnitude of internal variability (principally the frequency of ENSO (Kelley et al., 2020)) but broadly capture the observed trends. The response to volcanoes in 1982 and 1991 is clearer with the f2 and f3 forcing, demonstrating the impact of the stratospheric ozone chemistry correction from f1 even on surface temperatures. The E2.1-H model (same forcings, but with a different ocean model) and the E2.2-G model (with higher vertical resolution) have slightly lower (and more realistic) estimates of internal variability. Note



Figure 2. Tropospheric trends in various configurations showing the SAT, TLT and TMT changes, specifically the E2.1 AMIP; E2.1-G f1, f2, and f3; E2.1-H f2; and E2.2-G configurations for 1979–2014 (or to 2021 where available). The three diagnostics are offset by 1°C for clarity. The spread is the 95% confidence interval of the envelope of individual ensemble members. Observations from GISTEMP, UAH, RSS and NOAA STAR are in solid, dashed, dotted and dash-dotted lines, respectively.

		Ensemble	e Mean Linea	ar Trends ($^{\circ}C$	(dec)		
		1979-2014		,	1979 - 2021		
Configuration	SAT	TIT	TMT	SAT	TLT	TMT	
E2.1 AMIP f2	0.15^*	0.18^R	$0.12^{R,S}$				
E2.1-G fl	0.22^G	0.24^R	0.20^S				
E2.1-G f2	0.22^{*}	0.23^R	$0.17^{R,S}$	0.22^{*}	0.25^{R}	0.21	
E2.1-H f2	0.24	0.25^{R}	0.19^S	0.26^E	0.27^R	0.23	
E2.1-G f3	0.18^{*}	0.20^R	0.14^{*}				
E2.1-G (OMA)	0.18^{*}	0.20^R	$0.14^{R,S}$	0.19^*	0.22^R	$0.18^{R,S}$	
E2.1-G (MATRIX)	0.22	0.25^{R}	0.18^S	0.24	0.27	0.23	
E2.2-G	0.18^{*}	0.18^{*}	0.11^{*}				
E2.2-G (OMA)	0.14^{*}	0.15^*	0.07*	0.17^{*}	0.17^{*}	0.12^{*}	
Observations							
ERA5	0.16 ± 0.03			0.19 ± 0.03			
GISTEMP	0.16 ± 0.03			0.19 ± 0.02			
UAH v6		$0.11 {\pm} 0.05$	$0.07{\pm}0.05$		$0.13 {\pm} 0.04$	$0.10{\pm}0.04$	
RSS v4		0.19 ± 0.05	0.11 ± 0.05		$0.21{\pm}0.03$	$0.14{\pm}0.04$	
NOAA STAR v4.1			0.14 ± 0.05			$0.16 {\pm} 0.04$	

provide distributions from which any of the observations might plausibly be drawn (using Eqn. 1) are noted with a *. Where there is a difference depending on the trend is the 95% confidence interval on the linear regression. Trends through to 2021 (using the SSP2-4.5 simulations as a continuation) are available for some con-**Table 3.** Model and observed trends (^oC/dec) in tropospheric diagnostics. Model trends are derived from the ensemble mean. Uncertainty in the observational figurations. Ensemble mean trends that are consistent with at least one observational product within the observational uncertainty are in bold. Ensembles that observational product, the consistent product(s) is/are noted (G for GISS, E for ERA5, U for UAH, R for RSS, S for NOAA STAR).

$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	9-2014 $1979-2021$
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	TMT/SAT TLT/SAT TMT/SAT
$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$] 0.83 [0.76,0.88]
$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$	$0.92 \ [0.82, 1.01]$
$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$	[] 0.81 [0.71, 0.89] 1.11 [1.06, 1.18] 0.93 [0.88, 1.00]
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	$ \begin{bmatrix} 0.80 & [0.71, 0.87] \\ 0.80 & [0.71, 0.87] \end{bmatrix} \begin{bmatrix} 1.04 & [0.98, 1.10] \\ 0.69 & [0.56, 0.75] \end{bmatrix} $
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	0.73 [0.63, 0.80]
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	0.76[0.61, 0.87] 1.13 $[1.05, 1.20]$ 0.85 $[0.74, 0.98]$
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	0.83[0.76,0.89] 1.14 $[1.06,1.21]$ 0.95 $[0.86,1.02]$
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	0.59 [0.41, 0.78]
	0.45[0.21,0.61] 1.03[0.97,1.07] 0.69[0.56,0.75]
$\begin{array}{ c c c c c c c c c c c c c c c c c c c$	
$\begin{array}{ c c c c c c c c c c c c c c c c c c c$	
$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$	0.55 [0.41, 0.70] 0.74 [0.66, 0.81] 0.55 [0.47, 0.64]
$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	0.53 $[0.38, 0.69]$ 0.75 $[0.66, 0.83]$ 0.56 $[0.46, 0.65]$
GISTEMP+RSS 1.23 [1.13,1.32] 0.81 [0.67,0.94] 1.16 [1 ERA5+NOAA STAR 0.95 [0.84,1.06]	0.82 [0.70, 0.95] $1.13 [1.08, 1.19]$ $0.77 [0.69, 0.84]$
ERA5+NOAA STAR 0.95 [0.84,1.06]	[] 0.81 [0.67, 0.94] 1.16 [1.10, 1.21] 0.78 [0.69, 0.86]
	$0.95 \ [0.84, 1.06] \ 0.93 \ [0.79, 0.93]$
GISTEMP+NOAA STAR 0.94 [0.81,1.06]	0.94 [0.81, 1.06] 0.87 $[0.79, 0.95]$

s, with the 95%envelope from the ensemble. Uncertainty in the observational ratios is the 95% confidence interval on the linear regression through the origin of the two annual timeseries. Trend ratios through to 2021 (using the SSP2-4.5 simulations as a continuation) are available for some configurations. **Table 4.** Model and observ

that in the coupled models, the timing of ENSO variability will not in general be correlated with the observations.

Quantitative comparisons of the trends (in the historical period 1979–2014, and also in the extended period to 2021 for those ensembles that were continued under SSP2-4.5) can be seen in Table 3 for all nine model configurations. These results demonstrate more finely that there are notable differences in the trends among the configurations (and also the observational products) even within a broad qualitative agreement. For reference, differences in the ensemble mean trends that are greater than about $0.02\pm$ are statistically significant.

Before we address why specific ensembles have different trends, it's worth noting that even for a multi-decadal trend in the global mean temperature, there is significant spread across the individual ensemble members. For E2.1-G f2, for which we have ten simulations, the 1979–2021 SAT trends range from 0.18°C/dec to 0.26°C/dec, so even with 43 years of data, the spread can be important (see also Fyfe et al. (2021)). The analogous range for the E2.1-H and E2.2-G (OMA) configurations are [0.22,0.30] and [0.15,0.18]°C/dec respectively.

In comparing the model ensemble to the real world signal, the appropriate consis-334 tency test is whether the real world trend is a plausible draw from the modelled distri-335 bution. Thus even if the trend in the ensemble mean is outside the confidence interval 336 for the observed trend, the real world changes are potentially still consistent with the 337 modeled distribution (Santer et al., 2008). Bolding in Table 3 is based on whether the 338 ensemble mean is within the uncertainty of the observed trend (a test of whether the ob-339 served trend is consistent with our estimate of the forced trend), while the superscripts 340 denote whether each observational product can be considered a plausible draw from the 341 specific modeled ensemble. 342

Most configurations have a reasonable ensemble mean estimate of the tropospheric trends, and are consistent with ERA5, GISTEMP SAT and RSS TLT products. Only two configurations are also consistent with the UAH TLT estimates (both versions of E2.2-G model). Given the wide spread in estimated TMT trends across the observations, this diagnostic is less discriminating, though notably again, the E2.2-G configurations are consistent with the UAH trend.

Three configurations have ensemble SAT trends significantly greater than observations (and for two of them, this is also true for the longer 1979–2021 trend); E2.1-G f2, E2.1-H f2, and the E2.1-G (MATRIX) configurations. For E2.1-H, the Transient Climate Response (TCR) is higher than for the other configurations (Table 2) due to a reduced uptake of heat into the ocean (compared to the E2.1-G configurations), while the E2.1 (MATRIX) simulations have a more rapid decrease of (negative) aerosol forcings than with the OMA and non-interactive versions (Bauer et al., 2020).

Figure 3 show the 1979–2014 trends for the configurations. For each ensemble we show the density plot for each ensemble, with the exception of E2.1 (AMIP) which has a very narrow distribution. Across the configurations there is a wide range of trends for each diagnostic, skewed slightly higher than the observations, but on the whole mostly consistent with the RSS and NOAA STAR products. There are significant differences in the spread for specific configurations.

To better assess reasons for the spread in the trends, it's useful to also calculate the ratios of trends in the troposphere which removes the issues of overall global forcing and global mean temperature response, and allows for a focus on the global mean lapse rate, which is possibly more sensitive to model convective processes and other parameterizations (Santer et al., 2017). Table 4 gives the TLT to SAT, and TMT to SAT, ratios for each ensemble for the historical period and for the extended period to 2021. The ratios of the trends in the troposphere are relatively stable and similar to those seen in the CMIP5 multi-model ensemble (Santer et al., 2017). Whether these ratios in the models are calculated using the ensemble mean, or the mean of the individual trends for each ensemble member, or from a linear regression passing through the origin and through the individual annual points from each ensemble member, the values are basically the same. There is some suggestion that the trends get larger over time as the climate change signal increasingly dominates over the internal variability.

The observation-based trend ratios depend strongly on whether we use the UAH 375 data or the RSS/NOAA STAR data. With respect to UAH, the TLT/SAT trend ratio 376 is around 0.75, which is contrary to our basic physical understanding of the lapse rate, 377 indicating that there is likely a systematic problem in either all the SAT products or specif-378 ically the UAH TLT product. With respect to the other data products, the ratio is around 379 1.2, significantly greater than one (as expected). The TMT trends (and hence ratios) are 380 smaller because of the greater influence of the (cooling) stratosphere, but vary widely 381 across the products from around 0.55 (using UAH), 0.8 (using RSS) or 0.9 (using NOAA 382 STAR). 383

In the models, the TLT to SAT ratios are all greater than one, slightly less than 384 observational trend ratio using RSS, but entirely inconsistent with the trend ratio us-385 ing UAH. There is little spread in this value across different time periods, model struc-386 ture or forcing. However, there is a much wider spread in the trend ratios for the TMT 387 to SAT ratio, which vary by a factor of two between E2.1-G f1 and E2.2-G (OMA). This 388 result is tied to the spread in the TLS trends (see below), with the model configurations 389 with the largest cooling trends in the lower stratosphere having the smallest TMT/SAT 390 trend ratios. This underscores the importance of ozone, volcanic aerosols and possibly 391 even solar forcing in affecting the TMT trends and the TMT/SAT trend ratio. The E2.1-392 G f1 configuration had an error in the stratospheric ozone chemistry and the smallest 393 TLS cooling (to 1998), the f2 runs were corrected and had a more realistic stratospheric 394 cooling and a TMT/SAT trend close to that seen with the RSS products. 395

E2.1-G f2 visibly shows much clearer agreement with historical observations than E2.1-G f1. While the f2 improvements are easily noticeable in the period from the 1980's through the 90's, there is also a notable discrepancy in predicted and observed warming in the late 2000's and early 2010's.

The contrast between the E2.1-G f2 and f3 configurations is also instructive. These runs differ only in the aerosol and ozone fields, and show that the trends are quite sensitive to plausible changes in the aerosol forcing in particular. The TMT/SAT trend ratios however show more difference even though the underlying model processes are identical. This is plausibly connected with a lower tropopause in the E2.2 simulations from which the f3 forcings are drawn, implying a greater stratospheric contribution to the TMT trends.

The E2.1-G f2 ensemble runs hindcast more warming in the 2000's than the AMIP ensemble does. As the AMIP ensemble mean in use in these figures is the same as E2.1-G f2 runs except for its reliance on observational surface ocean temperature data, this indicates that, at least in part, the atmospheric-ocean dynamics in the coupled ensemble are contributing to more warming. However, it should be noted that the 95% confidence envelope of E2.1-G f2 does still overlap with the satellite observations in the majority of the temperature anomaly charts even as the model approaches the present.



Surface Air Temperature Decadal Trends 1979–2014

Figure 3. Trend analysis across the non-interactive configurations for the tropospheric diagnostics for the period 1979–2014. Uncertainties on the observational trends are the 95% confidence intervals. E2.1 (AMIP) results are plotted as a 95% spread. All other model ensembles are plotted as a density plot.



Figure 4. Stratospheric time-series for MSU TLS and SSU products for E2.1-G f1, f2 and f3, and for E2.2-G compared to the observations (each offset by 1°C for clarity). Observations from RSS, UAH and NOAA STAR are dotted, dashed and solid respectively.

										0.10 vith lse e
2411-3	0-000		-0.87	-0.88		-0.62^{*}	-0.57*		-0.59^{*}	-0.56± -1998 (w od, we u rvationi *. Wher
20	7-0		75	92		55*	47*		46^{*}	50 ± 0.10 nd post- led peric the obse with a $*$
2021/20			0-	-0.		-0-	-		0	38 -0.5 5 pre- a 5 extence within 7 noted STAR).
1999	1-000		-0.60	-0.61		-0.39^{*}	-0.30^{*}		-0.18^{*}	-0.40±0. t period . For the product rrawn are NOAA
SIL	CULT		-0.02*	-0.03*		$0.10^{R,S}$	0.14	-	0.43	$\left \begin{array}{c} -0.10\pm0.06\\ -0.05\pm0.06\\ -0.06\pm0.06\\ \end{array} \right $ pr two distinc ear regression observational blausibly be c or RSS, S for
SGI1_3	6-000	-0.77 -0 74	-0.77	-0.79	-0.61^{*}	-0.69*	-0.65^{*}	-0.57*	-0.53*	-0.51±0.15 d are given fi val in the lind at least one c tions might _F or UAH, <i>R</i> fi
ls (°C/dec) -2014 SSTL-2	7-000	-0.82 -0 77	-0.82	-0.84	-0.52^{*}	-0.65	-0.58^{*}	-0.49^{*}	-0.41*	-0.44 \pm 0.16 in and spread idence interv istent with ϵ the observat ϵ noted (U fi
Linear trend 1999- SSII-1	1-000	-0.76 -0.68	-0.75	-0.76	-0.33^{*}	-0.53^{*}	-0.44*	-0.31^{*}	-0.14^{*}	-0.37±0.13 asemble mea he 95% conf hat are cons which any of luct(s) is/are
S,		-0.25	$-0.20^{U,S}$	$-0.20^{U,S}$	0.23	-0.03*	0.04^*	0.23	0.34^{*}	$\left \begin{array}{c} -0.09\pm0.11\\ -0.04\pm0.11\\ -0.07\pm0.12\\ \text{agnostics. E}_{1}\\ \text{agnostics. E}_{2}\\ \text{tal trends t}\\ \text{thends t}\\ \text{thions from v}\\ \text{vsistent proc}\\ \text{asistent proc}\\ \end{array}\right.$
SGI1_2	6-000	-1.60 -1.58	-1.60	-1.63	-1.55	-1.65	-1.66	-1.51	-1.54	-1.11±0.29 tospheric di observatior mels. Ensen vide distribu duct, the col
1998 SSIL-2	7-000	-1.28 -1.1 3 *	-1.28	-1.31	-1.14^{*}	-1.34	-1.35	-1.15^{*}	-1.26^{*}	-1.02±0.23 dec) in stra tainty in the he SSU chan he SSU chan ples that pro vational prov
1979– SSIL-1	1-000	-0.90 * -0.64	-0.90*	-0.92^{*}	-0.76*	-0.95^{*}	-0.98*	-0.80*	-1.01^{*}	-0.86±0.22 trends (°C/ 021). Uncer d 2020 for t old. Ensemb n the obser
S IL	CULT	-0.25* -0.10	-0.28*	-0.26^{*}	-0.56^{*}	-0.36^{*}	-0.42^{*}	-0.61*	-0.85^{U}	-0.49±0.27 -0.41±0.28 -0.46±0.28 .nd observed 14 or 2020/2 .LS data, an llighted in bo depending o
Configuration	Comgutation	E2.1 AMIP f2 E2.1-G f1	E2.1-G f2	E2.1-H f2	E2.1-G f3	E2.1-G (OMA)	E2.1-G (MATRIX)	E2.2-G	E2.2-G (OMA)	Observations UAH v6 RSS v4 NOAA STAR Table 5. Model a trends ending in 20 up to 2021 for the 7 uncertainty are high there is a difference

-17-



Figure 5. Ensemble trends and observations for the TLS and SSU channels in the stratosphere. Uncertainties in the observations are the 95% confidence interval on the linear regression.

414 5.2 Stratospheric Trends

Figure 4 shows a selection of model anomalies in the stratospheric diagnostics, TLS 415 and the three SSU channels. The effect of the ocean module is not significant, so the re-416 sults of the E2.1-H or E2.1 (AMIP) configurations are omitted for clarity. Overall cool-417 ing trends increase as a function of height as the CO_2 impact increases, while the im-418 pact of volcanic aerosols decreases. It's clear that a single linear trend is not a good fit 419 for these diagnostics over this period because of the impacts of volcanic aerosols early 420 in the period, the changing impact of ozone depletion (strong in the first 20 years of the 421 record, and less important subsequently, see fig. 1) and the impact of solar cycles. The 422 overall structure of the changes is reasonable in all configurations, but there are clear 423 discrepancies. The difference in the volcanic response between the E2.1-G f1 and f2 sim-424 ulations is directly related to the correction of bug in the stratospheric chemistry, which 425 clearly improved the simulations. However, in all cases, the response to El Chichon in 426 1982 is muted compared to that for Mt. Pinatubo (1991) which appears to be too large. 427 The other clear difference is the between the f2 and f3 runs which have markedly dif-428 ferent upper stratospheric trends. The high-top simulations using E2.2-G have a much 429 closer match to the observations than E2.1-G f2, which is mimicked by the results in E2.1-430 G f3 which uses the E2.2-G (OMA) ozone fields, but has the same radiative effects from 431 CO_2 . 432

More quantitatively, the ensemble mean linear trends in the stratosphere are given 433 in Table 5. The ozone depletion signal is dominant in the TLS trends so, following Seidel 434 et al. (2016), we separate the stratospheric linear analysis into two periods, 1979–1998 435 (the 'ozone depletion' period) and the subsequent evolution 1999–2014 (or with a con-436 tinuation to 2021 or 2022). In the lower stratosphere, the different E2.1-G configurations 437 behave similarly in the early period. However, trends in the recovery period vary as a 438 function of the specific forcings. The f2 forcings indicate more cooling than other con-439 figurations in the SSU channels. This is not the case in the lower stratosphere, though, 440 which is driven more by the ozone trends. E2.1-G f2 results track the TLS observations 441 more closely than E2.1-G f3 or E2.2 do. Notably, in the decadal trend diagrams (fig. 5), 442 the E2.1-G f3 configuration has the widest difference between the ozone depletion and 443 recovery period while f2 is more centered around the observational data without signif-444 icant differentiation between the two periods. The improvements in model agreement 445 from the E2.1-G f2/f3 to E2.2 models configurations is increasingly pronounced in the 446 higher altitudes. 447

While the internal variability in the stratospheric trends is greatly muted compared with the troposphere, there is significant spread in each ensemble for the two periods (fig. 5) due to the different forcings. With the exception of SSU-3 (the upper stratospheric channel) in the ozone depletion period, the model ensembles using the high-top model (or the ozone fields derived from it) are a much better match to the observed trends. All model cooling trends in the ozone depletion period in SSU-3 are too large.

Some insight into the reasons behind the trend disparities can be gained by look-454 ing at the ozone trends in the different ensembles (fig. 6). We look at the $60^{\circ}S-60^{\circ}N$ mean 455 total column ozone (since the satellite observations don't completely capture the changes 456 at the high latitudes). The model ensembles are slightly depleted total ozone in the early 457 1980's (by 7-9 DU) compared to the SBUV (v8.7) data, but all show a steady decline 458 over the 'ozone depletion' period in line with the observations at least through to 1994. 459 In the E2.1 configurations, the depletion period lasts longer than observed (by 5 years 460 or so), while in the E2.2 configuration there is a greater perturbation associated with Mt 461 462 Pinatubo and a deeper depletion towards the the end of the 1990's.

In these ModelE simulations, stratospheric chlorine loading (and hence overall ozone
depletion) was set using a relationship based on the concentrations of CFC-11 and CFC12. In the real world, the Equivalent Effective Stratospheric Chlorine (EESC) (Newman



Figure 6. Time-series of the 12-month running mean total column ozone (DU) (60°S-60°N) for the ensembles with interactive composition from 1979–2021. Thick lines are the ensemble mean, with individual members in the thin lines. The vertical dashed line denotes the separation of the 'ozone depletion' and 'ozone recovery' periods. Observations are from SBUV v8.7 (McPeters et al., 2013, and updates).

Model Configuration	Intercept	Linear trend (°C/dec)	Volcanic AOD	Solar TSI
E2.1-G f2	-460	-1.15	4.4	0.48
E2.1-G f3	-563	-0.98	5.6	0.56
E2.2-G	-463	-1.02	5.1	0.49
E2.2-G (OMA)	-445	-1.13	5.8	0.49
NOAA STAR	-549 ± 240	-0.87 ± 0.11	$4.9{\pm}1.9$	$0.53 {\pm} 0.17$

SSU-2 Multiple Linear Regression Coefficients 1979–1998

Table 6. Multiple linear regression results for the SSU-2 diagnostics for selected ensemble mean model configurations and observations for the 'ozone depletion' period. All coefficients are significant at the 95% level. Uncertainties on the regression of the NOAA STAR observations are the 95% confidence levels.

et al., 2007) depends on many lower concentration gases which are not explicitly tracked 466 in the GCM. While the parameterized EESC is a good approximation to about 2000, 467 there is an increasing divergence with the real world afterwards, with the real EESC reducing more rapidly than in the model. Notably, the peak in the real world was around 469 2001, while in our parameterization it does not occur until 2005, and the real world EESC 470 has decreased by 14% from its peak in 2020, while in our parameterization, it has only 471 decreased 8%. Thus in all the model configurations, the ozone depletion period extends 472 a few more years than observed, and does not recover as quickly. Given the cooling im-473 pact of ozone depletion on the TLS and SSU channels, this issue can explain a portion 474 of the mismatched trends in the ozone recovery period. 475

476

5.2.1 Multiple linear regression using volcanic and solar predictors

The clear impacts of volcanic eruptions and solar cycles in the stratospheric diagnostics increase the non-linearity of the temperature trends in the stratosphere. We therefore redo the linear regressions including predictors for these effects to assess whether the longer term trends are being affected by potential errors in the modeling of the volcanic or solar responses. If those errors are significant, we should see a better match in the modeled and observed linear trends.

We use a volcanic predictor based on the aerosol optical depth history (Sato et al., 483 1993, and updates) and a solar index derived from a historical total solar irradiance dataset 484 (Coddington et al., 2016). We highlight the results with respect to SSU-2 in figure 7 and 485 for the ensemble means for all diagnostics in figure 8. As expected, the linear trends for 486 both the models and the observations in the ozone depletion period are both smaller in 487 magnitude and less uncertain when the volcanic and solar predictors are included. The 488 added predictors make the most difference in the TLS trends where the resulting linear 489 trends are all more consistent. 490

For the SSU channels the impacts are more muted, but there aren't any major shifts in the consistency with the observations. Notably, the E2.1-G **f2** simulations are show notably stronger stratospheric cooling than observed regardless of which additional predictors are included. The E2.1-G **f3** and E2.2-G results for the SSU-3 channel get closer to the observed trends, but are still too strong, suggesting that further investigation of the time-series of ozone depletion is required.

There are some interesting aspects of these regressions in the ozone depletion period. This is illustrated for SSU Channel 2 in Table 6. First, the solar regressions (ef-





Figure 7. Ensemble linear trends with and without volcanic aerosol and solar predictors in the mid-stratosphere (SSU-2) (1979–1998). Uncertainties on the observations represent a 95% confidence interval on the linear regression. a) The linear trends in each ensemble in the multiple linear regression and, b) trends in the ensemble mean for each configuration with and without the MLR.

499	fectively over two solar cycles are in line with the inference from the observations, as are
500	the volcanic effects. It's noticeable that the volcanic signal is stronger in the E2.2-G (OMA)
501	and E2.1-G f3 configurations than in the other two configurations and the observations
502	(though all coefficients are consistent with the observations). In all cases the linear trends
503	are now more coherent with the observations, but collectively they are still a little too
504	steep (except for E2.1-G f3 which is just about compatible). The slightly enhanced so-
505	lar effects in E2.1-G f3 may arise from the lack of solar-cycle related changes in photol-
506	ysis which would causes upper stratospheric water vapor to decrease at solar maxima,
507	damping the temperature impact.



Figure 8. Ensemble linear trends with and without using volcanic and solar predictors for the stratospheric diagnostics (1979–2014). Error bars on the observations represent a 95% confidence interval on the linear regression.



Impact of individual drivers in the E2.1-G f2 ensemble

Figure 9. Breakdown of the ensemble mean SAT, TMT, TLS, and SSU-2 anomalies for E2.1-G f2 as a function of relevant single forcings from 1880–2021 with respect to a baseline of 1880–1910. The uncertainty on the 'All drivers' line is the derived from the 95% confidence interval from the pre-industrial control run, which in practice is indistinguishable from the envelope of the individual ensemble member spread. For the tropospheric diagnostics, we apply a 4-year running mean filter to reduce the 'weather' noise that still remains in the ensemble mean for each single forcing (which only used 5 ensemble members). Illustrative observations are the GISTEMP LOTI, RSS MSU, and the NOAA-STAR SSU. SAT observations are plotted with the same baseline as the models, but for the satellite era diagnostics we align them so that their mean is equal to the model 'All drivers' ensemble mean over 1980–1999.

508 6 Single-forcing results

For the E2.1-G f2 (NINT) configuration, we performed a complete set of single forc-509 ing simulations for the historical period (5 ensemble members each). Thus for each di-510 agnostic, we can illustrate the modeled response to each of the drivers individually (fig. 9). 511 Because of the way the historical composition files were derived (from an OMA simu-512 lation), the solar-only and volcanic-only forcing simulations include compositional responses 513 (notably in ozone) which is a new feature compared to how similar exercises were done 514 in previous iterations (Marvel et al., 2015). Some forcings (such as orbital forcing or land-515 use/land-cover change) don't have a significant expression in the global mean diagnos-516 tics (very close to zero for orbital forcing, and slightly negative for land-use/land-cover 517 on SAT) and are not included in the figures. The impact of 'Tropospheric Aerosols' is 518 only significant for the tropospheric diagnostics, though there is a very slight stratospheric 519 warming associated with them (not shown). The ozone-only results used composition 520 files from anthropogenic-only simulations (with no natural drivers), and thus include both 521 tropospheric ozone increases and stratospheric ozone decreases, driven by emissions of 522 chemical precursors and ozone-depleting substances. The Greenhouse Gas (GHG)-only 523

simulations include the radiative impacts of CO_2 , CH_4 , N_2O and CFCs (but not any chemistry related impacts).

For the SAT, TMT and SSU2 diagnostics illustrated in figure 9, the dominance of 526 GHGs in driving the long-term trend is clear, however, other forcings (tropospheric aerosols, 527 ozone, volcanic aerosols) all play key roles, though their importance varies through the 528 atmosphere. Volcanic and ozone forcings are relevant throughout the atmosphere, while 529 solar forcing increases in importance with height. GHG, volcanic and ozone impacts all 530 change sign in going from the troposphere to the stratosphere. The breakdown for TLT 531 532 is similar to that for SAT, and the two other SSU channels resemble SSU-2 (not shown). For TLS, the ozone changes are the dominant factor in recent decades (c.f. Ramaswamy 533 et al. (1996), although the GHG (CO_2) impact is increasing in importance. The vari-534 ations across the other model configurations, particularly in the stratosphere, can be thought 535 of as being driven by small changes in the secondary components - notably stratospheric 536 ozone and the volcanic response. 537

538 7 Discussion

The vertical profile of atmospheric trends over recent decades is a key metric in assessing the fidelity of climate models, and ultimately in understanding why the current climate is changing. While the overall patterns are robust and clear - warming in the troposphere, cooling in the stratosphere, punctuated by volcanic effects, and modulated by solar activity - there are sufficient discrepancies between models and observations and among observational products to merit closer attention.

Among the dozens of simulations with the GISS Earth System Models in nine different configurations, there is sufficient structural variety to help us more easily identify some key modeling choices that have impacted those comparisons than when looking across the whole multi-model ensemble.

Most obviously, even for trends over four decades, there is substantial intra-ensemble spread due to the different realizations of internal variability in the troposphere and which is essential to take into account when comparing a model to the single real world realization (Santer et al., 2008; Po-Chedley et al., 2021).

Secondly, two factors for which there is still substantial uncertainty - the tropospheric 553 aerosol forcing changes and the rate and manner of heat uptake into the ocean - still make 554 a notable difference in the troposphere. Model configurations with less deep ocean heat 555 uptake and those with a faster decrease in aerosols have stronger surface trends than those 556 without. Additionally, while the spread of climate sensitivity in these configurations (2.1°C 557 to 3.1° C) is not as wide as in the broader CMIP6 ensemble ($1.8-5.6^{\circ}$ C) (Zelinka et al., 558 2020), it is sufficient for the trends under similar forcings to diverge. Unfortunately, the 559 intersection of these three factors means that it is hard to constrain one of them alone. 560

Thirdly, it is likely that there will be further refinements and better estimates of structural uncertainty for the satellite retrievals themselves. If so, the conclusions here may need to be revised.

Nonetheless, there are robust conclusions that can be drawn. There is clearly more 564 work to be done related to the response of the models to volcanic eruptions. The dis-565 crepancies seen in the magnitude and time-evolution of volcanic signal suggest that ei-566 ther the input aerosol fields are not accurate, and/or that the model response (perhaps 567 in the heterogeneous chemistry) is flawed. More first principles modeling via injection 568 of volcanic gases and subsequent aerosol modeling (LeGrande et al., 2016) and the re-569 sults of the VolMIP exercise might lead to more coherent and hopefully more accurate 570 impacts (Timmreck et al., 2018; Zanchettin et al., 2022). 571

Our results underline the importance of ozone climatology and chemistry responses. 572 The difference in SSU trends in the E2.1-G f2 and f3 configurations can only be due 573 to the different ozone files. There are two facets to these differences, a more accurate base 574 climatology, with a lower altitude ozone layer (consistent with a more accurate (weaker) 575 Brewer-Dobson circulation and older stratospheric age-of-air in the E2.2 models (Orbe 576 et al., 2020)) and different trends over time. We find that E2.2-G is much better than 577 E2.1-G in its agreement with ozone depletion and recovery observations. Our results sug-578 gest that the magnitude of the SSU cooling trends (driven mainly by the CO₂ forcing) 579 are mediated by the ozone response in the models. Ozone responds differently in the high-580 top versus low-top models because of the climate changes in the Brewer-Dobson circu-581 lation and because the slower circulation in the high-top versions changes the response 582 of ozone to increased CO_2 . Improvements in the modeling of stratospheric halogen loads 583 may also make a difference to the trends in the 'ozone recovery' period. 584

It should be noted that, in agreement with other CMIP6 models, GISS ModelE out-585 put tends to agree more with RSS and NOAA STAR tropospheric observations than with 586 the UAH data. The results from the AMIP results are very suggestive that the UAH re-587 sults start to anomalously deviate from expectations around the year 2000. Updates to 588 the SST inputs for the AMIP simulations are likely to worsen the comparison further. 589 Also, since the ratio of TLT to SAT trends is a robust metric across configurations, in-590 dependent of the climate sensitivity, vertical resolution or ocean component, the fact that 591 this is not consistent with any UAH/SAT trend ratio is suggestive of a systematic prob-592 lem. 593

To summarise, it is too simplistic to attribute all model discrepancies with the MSU 594 and SSU observational to a single dominant cause. This analysis has demonstrated that 595 even within a single model family, multiple factors are at play: ozone chemistry, clima-596 tology and feedbacks are clearly important to the TMT and stratospheric channels; cor-597 rect simulations of volcanic aerosol and consequent compositional changes are also im-598 portant. Both are targets for future development. However, internal variability and struc-599 tural uncertainty in the observations are essential components to address in any anal-600 ysis. Attempts to classify the responses in the multi-model ensemble by using only sin-601 gle ensemble members from each model or model family will, simply by chance, conflate 602 internal variability with structural uncertainty (e.g. McKitrick and Christy (2020); Mitchell 603 et al. (2020)) and may give misleading results. 604

605 Open Science

The observed MSU/SSU data products are available at ftp://ftp.remss.com/ 606 msu/graphics/, https://www.nsstc.uah.edu/data/msu/v6.0/ and ftp://ftp.star 607 .nesdis.noaa.gov/pub/smcd/emb/mscat/data/. The GISTEMP data is available from 608 https://data.giss.nasa.gov/gistemp. ERA5 data is from https://climate.copernicus 609 .eu/sites/default/files/ftp-data/temperature/2021/12/ERA5_1991-2020/ts_1month 610 _anomaly_Global_ERA5_2T_202112_1991-2020_v01.csv. Indicies of various drivers are 611 from the following sites: ozone hole area from NASA Ozone Watch https://ozonewatch 612 .gsfc.nasa.gov/statistics/annual_data.txt, total column ozone https://acd-ext 613 .gsfc.nasa.gov/Data_services/merged/data/sbuv_v87_mod.int_lyr.70-20.za.r1 614 .txt, volcanic aerosol optical depth https://data.giss.nasa.gov/modelforce/strataer/ 615 tau.line_2012.12.txt, Total Solar Irradiance https://lasp.colorado.edu/lisird/ 616 latis/dap/nrl2_tsi_P1Y.csv?&time%3E=1610-01-01T00:00:00.0002&time%3C=2021 617 -12-31T00:00:00.000Z, and radiative forcing https://data.giss.nasa.gov/modelforce/ 618 Miller_et_al21/ERFs_SSP245_MillerFig10_2021.txt. CMIP6 data is available from 619 the Earth System Grid Federation (ESGF) https://esgf-node.llnl.gov/search/cmip6/ 620 or from the NASA Center for Climate Simulations (NCCS) https://portal.nccs.nasa 621 .gov/datashare/giss_cmip6/. Non-CMIP6 simulations and derived data (such as the 622

- model MSU and SSU diagnostics) are available from https://portal.nccs.nasa.gov/
 datashare/giss_cmip6/NON-CMIP/.
- 625 Acknowledgments
- ⁶²⁶ Climate modeling at NASA GISS is supported by the NASA Modeling Analysis and Pre-
- diction program, and resources supporting this work were provided by the NASA High-
- End Computing (HEC) Program through NCCS at Goddard Space Flight Center (GSFC).
- ⁶²⁹ M.C.C. was supported by the Climate Change Research Initiative (CCRI) at GISS funded
- ⁶³⁰ by the GSFC Office of STEM Engagement.

631 References

643

644

645

- Abdul-Razzak, H., & Ghan, S. J. (2000). A parameterization of aerosol activation: 2. Multiple aerosol types. Journal of Geophysical Research: Atmospheres, 105 (D5), 6837–6844. doi: 10.1029/1999jd901161
- Abdul-Razzak, H., Ghan, S. J., & Rivera-Carpio, C. (1998). A parameterization of
 aerosol activation: 1. Single aerosol type. Journal of Geophysical Research: At mospheres, 103(D6), 6123–6131. doi: 10.1029/97jd03735
- Bauer, S. E., Tsigaridis, K., Gao, Y. C., Faluvegi, G., Kelley, M., Lo, K. K., ... Wu,
 J. (2020). Historical (1850–2014) aerosol evolution and role on climate forcing
 using the GISS ModelE2.1 contribution to CMIP6. Journal of Advances in
 Modeling Earth Systems. doi: 10.1029/2019MS001978
- Bauer, S. E., Wright, D. L., Koch, D., Lewis, E. R., McGraw, R., Chang, L.-S., ...
 - Ruedy, R. (2008). MATRIX (Multiconfiguration Aerosol TRacker of mIXing state): an aerosol microphysical module for global atmospheric models. *Atmos. Chem. Phys.*, 8, 6003–6035.
- CCSP. (2006). Temperature trends in the lower atmosphere: Steps for understanding and reconciling differences. Asheville, NC, USA: National Oceanic and Atmospheric Administration, National Climatic Data Center. (Karl, T. R. et al., eds, 164 pp.)
- ⁶⁵⁰ Christy, J. R., & Spencer, R. W. (1995). Assessment of precision in temperatures
 ⁶⁵¹ from the microwave sounding units. *Climatic Change*, 30(1), 97–102. doi: 10
 ⁶⁵² .1007/bf01093227
- ⁶⁵³ Coddington, O., Lean, J. L., Pilewskie, P., Snow, M., & Lindholm, D. (2016). A so lar irradiance climate data record. Bulletin of the American Meteorological So *ciety*, 97(7), 1265–1282. doi: 10.1175/bams-d-14-00265.1
- ⁶⁵⁶ Flannaghan, T. J., Fueglistaler, S., Held, I. M., Po-Chedley, S., Wyman, B., & Zhao,
 M. (2014). Tropical temperature trends in Atmospheric General Circulation
 ⁶⁵⁸ Model simulations and the impact of uncertainties in observed SSTs. Journal
 ⁶⁵⁹ of Geophysical Research: Atmospheres, 119(23). doi: 10.1002/2014jd022365
- Fu, Q., Johanson, C. M., Warren, S. G., & Seidel, D. J. (2004). Contribution of stratospheric cooling to satellite-inferred tropospheric temperature trends. *Nature*, 429, 55–58.
- Fu, Q., Manabe, S., & Johanson, C. M. (2011). On the warming in the tropical upper troposphere: Models versus observations. *Geophysical Research Letters*, 38. doi: 10.1029/2011gl048101
- Fyfe, J. C., Kharin, V. V., Santer, B. D., Cole, J. N. S., & Gillett, N. P. (2021).
 Significant impact of forcing uncertainty in a large ensemble of climate model
 simulations. *Proceedings of the National Academy of Sciences*, 118(23). doi: 10.1073/pnas.2016549118
- Hansen, J., Wilson, H., Sato, M., Ruedy, R., Shah, K., & Hansen, E. (1995). Satel lite and surface temperature data at odds? *Climatic Change*, 30, 103–117. doi: 10.1007/BF01093228
- Hersbach, H., Bell, B., Berrisford, P., Hirahara, S., Horányi, A., Muñoz Sabater, J.,

674	Thépaut, JN. (2020). The ERA5 global reanalysis. Quarterly Journal of
675	the Royal Meteorological Society, 146(730), 1999–2049. doi: 10.1002/qj.3803
676	Jones, P. D., New, M., Parker, D. E., Martin, S., & Rigor, I. G. (1999). Surface
677	air temperature and its variations over the last 150 years. Revs. Geophys., 37,
678	173–199.
679	Jones, P. D., Osborn, T. J., Wigley, T. M. L., Kelly, P. M., & Santer, B. D. (1997).
680	Comparisons between the microwave sounding unit temperature record and
681	the surface temperature record from 1979 to 1996: Real differences or potential
682	discontinuities? Journal of Geophysical Research: Atmospheres, 102(D25),
683	30135–30145. doi: 10.1029/97jd02432
684	Kelley, M., Schmidt, G. A., Nazarenko, L. S., Bauer, S. E., Ruedy, R., Rus-
685	sell, G. L., Yao, MS. (2020). GISS-E2.1: Configurations and cli-
686	matology. Journal of Advances in Modeling Earth Sustems, 12. doi:
687	10.1029/2019 ms 002025
688	Kennedy, J. J., Bayner, N. A., Atkinson, C. P., & Killick, R. E. (2019). An en-
689	semble data set of sea surface temperature change from 1850. The Met Office
690	Hadley Centre HadSST 4 0 0 0 data set <i>Journal of Geophysical Research</i> :
691	Atmospheres 12/(14) 7719–7763 doi: 10.1029/2018id029867
602	Kramarova N A Nash E B Newman P A Bhartia P K McPeters B D
602	Bault D F Labow G I (2014) Measuring the Antarctic ozone hole
604	with the new Ozone Mapping and Profiler Suite (OMPS) Atmospheric Chem-
605	istry and Physics $1/(5)$ 2353–2361 doi: 10.5194/acp-14-2353-2014
695	LeCrande A N Tsigaridis K & Bayer S E (2016) Role of atmospheric chem-
690	istry in the climate impacts of stratospheric volcanic injections Nature Geo-
697	science $q(0)$ 652–655 doi: 10.1038/ngeo2771
698	Lonsson N I I Schmidt C A Hanson I F Monno M I Porsin A Buody
699	B & Zyzz D (2010) Improvements in the CISTEMP uncertainty model
700	Lowrnal of Coonducted Research: Atmospheres 19/ 6307-6326
701	10 1020 /2018;d020522
702	Manaha S fr Wethereld P T (1067) Thermal equilibrium of the atmosphere
703	with a given distribution of relative humidity <i>I</i> Atmos Sci. 21 241 250
704	Marval K Schmidt C A Miller P. L. & Nazaranka I. S. (2015. doc) Impli
705	estions for alimete consistivity from the response to individual foreinga. Nature
706	Climate Change 6(4) 386 380 doi: 10.1038/polimeto2888
707	M_{assen} Delmette V et el (Edg.) (2021) <i>Climete Change 0021</i> . The physical sei
708	Masson-Definitie, V., et al. (Eds.). (2021). Climate Change 2021: The physical sci-
709	the Interconcernmental Dance on Climate Change, Combridge University Proce
710	Margarela A. C. Davidal W. J. Stainer, A. K. Kamashka, A. V. Christer, J. Sam
711	Maycock, A. C., Randel, W. J., Steiner, A. K., Karpechko, A. Y., Unristy, J., Saun-
712	terra anothing theory of recent stratospheric derived and the mystery of recent stratospheric derived and the strates of the s
713	temperature trends. Geophysical research Letters, $40(18)$, $9919-9933$. doi: 10.1020/2018~1078035
714	10.1029/2010g1070030 M-Kitrich D $\theta_{\rm c}$ (heister L (2020) Denversion action relation of CMIDC transcendence
715	MCKItrick, R., & Unristy, J. (2020). Pervasive warming bias in UMP6 tropospieric
716	ayers. Larth and space science. doi: 10.1029/2020ea001281
717	McPeters, R. D., Bhartia, P. K., Haffner, D., Labow, G. J., & Flynn, L. (2013). The
718	version 8.0 SBUV ozone data record: An overview. Journal of Geophysical Re-
719	search: Atmospheres, $110(14)$, $8032-8039$. doi: $10.1002/\text{Jgrd.}50597$
720	Wears, U. A., Schabel, M., & Wentz, F. J. (2003). A reanalysis of the MSU Channel 2 true exclusion terms and L L CL = 16, 2000 acct.
721	2 tropospheric temperature record. J. Clim., 1b, 3050–3064.
722	Mears, C. A., Wentz, F. J., & Thorne, P. W. (2012). Assessing the value of Mi-
723	crowave Sounding Unit-radiosonde comparisons in ascertaining errors in
724	climate data records of tropospheric temperatures. Journal of Geophysical
725	Research: Atmospheres, 117. doi: $10.1029/2012jd017710$
726	Menon, S., Koch, D., Beig, G., Sahu, S., Fasullo, J., & Orlikowski, D. (2010). Black
727	carbon aerosols and the third polar ice cap. Atmos. Chem. Phys., 10, 4559–
728	4571.

729	Menon, S., & Rotstayn, L. (2006). The radiative influence of aerosol effects
730	on liquid-phase cumulus and stratiform clouds based on sensitivity stud-
731	ies with two climate models. Climate Dynamics, $27(4)$, $345-356$. doi:
732	10.1007/s00382-006-0139-3
733	Menon, S., Unger, N., Koch, D., Francis, J., Garrett, T., Sednev, I., Streets, D.
734	(2008). Aerosol climate effects and air quality impacts from 1980 to 2030.
735	Environ. Res. Lett., 3. doi: 10.1088/1748-9326/3/2/024004
736	Miller, R. L., Schmidt, G. A., Nazarenko, L., Bauer, S. E., Kellev, M., Ruedy, R.,
737	Yao, MS. (2021). CMIP6 historical simulations (1850-2014) with GISS
738	ModelE2.1. J. Adv. Model. Earth Syst., 13. doi: 10.1029/2019MS002034
739	Miller, R. L., Schmidt, G. A., Nazarenko, L. S., Tausnev, N., Ruedy, R., Kel-
740	lev, M., Zhang, J. (2014). CMIP5 historical simulations (1850–
741	2012) with GISS ModelE2. J. Adv. Model. Earth Syst., 6, 441–477. doi:
742	10.1002/2013MS000266
743	Mitchell, D. M., Lo, Y. T. E., Seviour, W. J. M., Haimberger, L., & Polvani, L. M.
744	(2020). The vertical profile of recent tropical temperature trends: Persistent
745	model biases in the context of internal variability. <i>Environmental Research</i>
746	<i>Letters</i> , 15. doi: 10.1088/1748-9326/ab9af7
747	Mitchell, D. M., Thorne, P. W., Stott, P. A., & Grav, L. J. (2013). Revisiting the
748	controversial issue of tropical tropospheric temperature trends. <i>Geophysical Re</i> -
749	search Letters, 40, 2801–2806. doi: 10.1002/grl.50465
750	Morice, C. P., Kennedy, J. J., Ravner, N. A., Winn, J. P., Hogan, E., Killick, R. E.,
751	Simpson, I. R. (2021). An updated assessment of near-surface temperature
752	change from 1850: The HadCRUT5 data set. Journal of Geophysical Research:
753	Atmospheres, 126. doi: 10.1029/2019jd032361
754	Nazarenko, L. S., Tausnev, N., Russell, G. L., Rind, D., Miller, R. L., Schmidt,
755	G. A., Yao, MS. (2022). Future climate change under SSP emis-
756	sion scenarios with GISS-E2.1. J. Adv. Model. Earth Syst doi: 10.1029/
757	2021ms002871
758	Newman, P. A., Daniel, J. S., Waugh, D. W., & Nash, E. R. (2007). A new formula-
759	tion of equivalent effective stratospheric chlorine (EESC). Atmospheric Chem-
760	istry and Physics, 7(17), 4537–4552. doi: 10.5194/acp-7-4537-2007
761	Orbe, C., Rind, D., Jonas, J., Nazarenko, L., Faluvegi, G., Murray, L. T.,
762	Schmidt, G. A. (2020). GISS model E2.2: A climate model optimized for
763	the middle atmosphere—2. Validation of large-scale transport and evaluation
764	of climate response. Journal of Geophysical Research: Atmospheres, 125(24).
765	doi: 10.1029/2020jd033151
766	Po-Chedley, S., & Fu, Q. (2012). Discrepancies in tropical upper tropospheric warm-
767	ing between atmospheric circulation models and satellites. Environmental Re-
768	search Letters, 7. doi: 10.1088/1748-9326/7/4/044018
769	Po-Chedley, S., Santer, B. D., Fueglistaler, S., Zelinka, M. D., Cameron-Smith,
770	P. J., Painter, J. F., & Fu, Q. (2021). Natural variability contributes to
771	model–satellite differences in tropical tropospheric warming. Proceedings of the
772	National Academy of Sciences, 118(13). doi: 10.1073/pnas.2020962118
773	Ramaswamy, V., Schwarzkopf, M. D., & Randel, W. J. (1996). Fingerprint of ozone
774	depletion in the spatial and temporal pattern of recent lower-stratospheric
775	cooling. Nature, $382(6592)$, 616–618. doi: 10.1038/382616a0
776	Randel, W. J., & Cobb, J. B. (1994). Coherent variations of monthly mean total
777	ozone and lower stratospheric temperature. Journal of Geophysical Research,
778	99(D3), 5433. doi: 10.1029/93jd03454
779	Rayner, N. A., Brohan, P., Parker, D. E., Folland, C. K., Kennedy, J. J., Vanicek,
780	M., Tett, S. F. B. (2006). Improved analyses of changes and uncertain-
781	ties in sea surface temperature measured in situ since the mid-nineteenth
782	century: The HadSST2 dataset. Journal of Climate, $19(3)$, $446-469$. doi:
783	10.1175/jcli3637.1

784	Richardson, M., Cowtan, K., & Millar, R. J. (2018). Global temperature definition
785	affects achievement of long-term climate goals. Environmental Research Let-
786	ters, $13(5)$, 054004. doi: 10.1088/1748-9326/aab305
787	Rind, D., Orbe, C., Jonas, J., Nazarenko, L., Zhou, T., Kelley, M., Schmidt,
788	G. A. (2020). GISS model E2.2: A climate model optimized for the middle
789	atmosphere. Model structure, climatology, variability and climate sensitivity.
790	Journal of Geophysical Research: Atmospheres. doi: 10.1029/2019jd032204
791	Santer, B. D., Bonfils, C., Painter, J. F., Zelinka, M. D., Mears, C., Solomon, S.,
792	Wentz, F. J. (2014). Volcanic contribution to decadal changes in tropospheric
793	temperature. Nature Geosci., $7(3)$, 185–189, doi: 10.1038/ngeo2098
704	Santer B D Hnilo I I Wigley T M L Boyle I S Doutriaux C Fiorino
794	M Taylor K E (1999) Uncertainties in observationally based esti-
795	mates of temperature change in the free atmosphere I Geophys Res 10/
797	6305–6333.
798	Santer, B. D., Po-Chedley, S., Mears, C., Fyfe, J. C., Gillett, N., Fu, Q., Zou, C
799	Z. (2021). Using climate model simulations to constrain observations. Journal
800	of Climate, 1–59. doi: 10.1175/jcli-d-20-0768.1
801	Santer, B. D., Solomon, S., Pallotta, G., Mears, C., Po-Chedley, S., Fu, Q., Bon-
802	fils, C. (2017). Comparing tropospheric warming in climate models and satel-
803	lite data. Journal of Climate, 30(1), 373–392. doi: 10.1175/jcli-d-16-0333.1
804	Santer, B. D., Taylor, K. E., Wigley, T. M. L., Johns, T. C., Jones, P. D., Karoly,
805	D. J., Tett, S. (1996). A search for human influences on the thermal struc-
806	ture of the atmosphere. Nature, $382(6586)$, $39-46$. doi: $10.1038/382039a0$
807	Santer, B. D., Thorne, P. W., Haimberger, L., Taylor, K. E., Wigley, T. M. L., Lan-
808	zante, J. R., Wentz, F. J. (2008). Consistency of modelled and observed
809	temperature trends in the tropical troposphere. International Journal of Cli-
810	matology, 28(13), 1703–1722. doi: 10.1002/joc.1756
811	Santer, B. D., Wigley, T. M. L., Mears, C., Wentz, F. J., Klein, S. A., Seidel, D. J.,
812	Schmidt, G. A. (2005). Amplification of surface temperature trends and
813	variability in the tropical atmosphere. Science, $309(5740)$, $1551-1556$. doi:
814	10.1126/science.1114867
815	Sato, M., Hansen, J. E., McCormick, M. P., & Pollack, J. B. (1993). Stratospheric
816	aerosol optical depths, 1850–1990. J. Geophys. Res., 98, 22,987–22,994.
817	Schmidt, G. A., Kelley, M., Nazarenko, L., Ruedy, R., Russell, G. L., Aleinov, I.,
818	Zhang, J. (2014). Configuration and assessment of the GISS ModelE2
819	contributions to the CMIP5 archive. J. Adv. Model. Earth Syst., 6, 141–184.
820	doi: 10.1002/2013MS000265
821	Schmidt, G. A., Ruedy, R., Hansen, J. E., Aleinov, I., Bell, N., Bauer, M.,
822	Yao, MS. (2006). Present-day atmospheric simulations using GISS Mod-
823	elE:Comparison to in situ, satellite, and reanalysis data. Journal of Climate,
824	19, 153–192. doi: 10.1175/jcli3612.1
825	Seidel, D. J., Li, J., Mears, C., Moradi, I., Nash, J., Randel, W. J., Zou, CZ.
826	(2016). Stratospheric temperature changes during the satellite era. Journal of
827	Geophysical Research: Atmospheres, 664–681. doi: 10.1002/2015jd024039
828	Shah, K. P., & Rind, D. (1995). Use of microwave brightness temperatures with a
829	general circulation model. J. Geophys. Res., 100, 13,841–13,874.
830	Sheather, S. J., & Jones, M. C. (1991). A reliable data-based bandwidth se-
831	lection method for kernel density estimation. Journal of the Royal Sta-
832	tistical Society: Series B (Methodological), 53(3), 683–690. doi: 10.1111/
833	j.2517-6161.1991.tb01857.x
834	Simmons, A., Hersbach, H., Munoz-Sabater, J., Nicolas, J., Vamborg, F., Berrisford,
835	P., Woollen, J. (2021). Low frequency variability and trends in surface
836	air temperature and humidity from ERA5 and other datasets. ECMWF. doi:
837	10.21957/LY5VBTBFD
838	Spencer, R. W., & Christy, J. R. (1990). Precise monitoring of global temperature

839	trends from satellites. Science, 247, 1558–1562.
840	Spencer, R. W., Christy, J. R., & Braswell, W. D. (2017). UAH version 6 global
841	satellite temperature products: Methodology and results. Asia-Pacific Journal
842	of Atmospheric Sciences, 53(1), 121–130. doi: 10.1007/s13143-017-0010-y
843	Taylor, K. E., Williamson, D., & Zwiers, F. (2000). The sea surface tempera-
844	ture and sea ice concentration boundary conditions for AMIP II simula-
845	tions. PCMDI Report 60, Program for Climate Model Diagnosis and In-
846	tercomparison, Lawrence Livermore National Laboratory. Retrieved from
847	https://pcmdi.llnl.gov/mips/amip/index.html
848	Thompson, D. W. J., Seidel, D. J., Randel, W. J., Zou, CZ., Butler, A., Mears, C.,
849	Lin, R. (2012). The mystery of recent stratospheric temperature trends.
850	<i>Nature</i> , 491, 692–697. doi: 10.1038/nature11579
851	Thorne, P. W., Brohan, P., Titchner, H. A., McCarthy, M. P., Sherwood, S. C., Pe-
852	terson, T. C., Kennedy, J. J. (2011). A quantification of uncertainties in
853	historical tropical tropospheric temperature trends from radiosondes. Journal
854	of Geophysical Research, $116.$ doi: $10.1029/2010$ jd015487
855	Timmreck, C., Mann, G. W., Aquila, V., Hommel, R., Lee, L. A., Schmidt, A.,
856	Weisenstein, D. (2018). The interactive stratospheric aerosol model intercom-
857	parison project (isa-mip): motivation and experimental design. Geoscientific
858	Model Development, $11(7)$, 2581–2608. doi: 10.5194/gmd-11-2581-2018
859	Vose, R. S., Huang, B., Yin, X., Arndt, D., Easterling, D. R., Lawrimore, J. H.,
860	Zhang, H. M. (2021). Implementing full spatial coverage in NOAA's
861	global temperature analysis. Geophysical Research Letters, 48. doi:
862	10.1029/2020g1090873
863	Wentz, F. J., & Schabel, M. (1998). Effects of orbital decay on satellite-derived
864	lower-tropospheric temperature trends. Nature, 394, 661–664. doi: 10.1038/
865	$\frac{29267}{10000}$
866	Wigley, I. M. L. (2006). Appendix A: Statistical issues regarding trends. In
867	1. R. Karl, S. J. Hassol, C. D. Miller, & W. L. Multay (Eds.), <i>Temperature</i>
868	trends in the lower atmosphere: Steps for understanding and reconcurring differ-
869	Change Descende Weshington DC
870	Zanahattin D. Timmraak C. Khadri M. Sahmidt A. Taahay M. Aha M.
871	Weierbach H (2022) Effects of foreing differences and initial conditions on
872	inter model agreement in the VelMIP vole pinetube full experiment. <i>Geoscien</i>
873	tific Model Development 15(5) 2265-2202 doi: 10.5104/gmd 15.2265.2022
874	Zelinka M D Myers T A McCov D T Po-Chedley S Caldwall P M Cappi
875	P Taylor K E (2020) Causes of higher climate sensitivity in CMIP6
870	models Geonhusical Research Letters 17 doi: 10.1020/2010g1085782
878	Zou C-Z & Qian H (2016) Stratospheric temperature climate data record from
879	204, 0. 2., a gran, in (2010). Stratospheric temperature emilate data record nom
	merged SSU and AMSU-A observations. Journal of Atmospheric and Oceanic