

Using regional scaling for temperature forecasts with the Stochastic Seasonal to Interannual Prediction System (StocSIPS).

Lenin Del Rio Amador (✉ delrio@physics.mcgill.ca)

McGill University Department of Physics <https://orcid.org/0000-0003-4043-472X>

Shaun Lovejoy

McGill University

Research Article

Keywords: Macroweather Regime, Power-law Correlations, Stochastic Model, Energy Balance Equation, Spatially Resolved Temperature Field, Internal Variability

Posted Date: March 17th, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-326161/v1>

License:  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Version of Record: A version of this preprint was published at Climate Dynamics on April 5th, 2021. See the published version at <https://doi.org/10.1007/s00382-021-05737-5>.

1 Using regional scaling for temperature forecasts with the Stochastic 2 Seasonal to Interannual Prediction System (StocSIPS).

3 Lenin Del Rio Amador¹, Shaun Lovejoy¹

4 ¹Physics, McGill University, 3600 University St., Montreal, Que. H3A 2T8, Canada

5 *Correspondence to:* Lenin Del Rio Amador (delrio@physics.mcgill.ca) ORCID iD: 0000-0003-4043-472X

6 **Abstract.** Over time scales between 10 days and 10-20 years – the macroweather regime – atmospheric fields, including the
7 temperature, respect statistical scale symmetries, such as power-law correlations, that imply the existence of a huge memory
8 in the system that can be exploited for long-term forecasts. The Stochastic Seasonal to Interannual Prediction System
9 (StocSIPS) is a stochastic model that exploits these symmetries to perform long-term forecasts. It models the temperature as
10 the high-frequency limit of the (fractional) energy balance equation (fractional Gaussian noise) which governs radiative
11 equilibrium processes when the relevant equilibrium relaxation processes are power law, rather than exponential. They are
12 obtained when the order of the relaxation equation is fractional rather than integer and they are solved as past value problems
13 rather than initial value problems. StocSIPS was first developed for monthly and seasonal forecast of globally averaged
14 temperature. In this paper, we extend it to the prediction of the spatially resolved temperature field by treating each grid point
15 as an independent time series. Compared to traditional global circulation models (GCMs), StocSIPS has the advantage of
16 forcing predictions to converge to the real-world climate. It extracts the internal variability (weather noise) directly from past
17 data and does not suffer from model drift. Here we apply StocSIPS to obtain monthly and seasonal predictions of the surface
18 temperature and show some preliminary comparison with multi-model ensemble (MME) GCM results. For one month lead
19 time, our simple stochastic model shows similar values of the skill scores than the much more complex deterministic models.

20 1 Introduction

21 The Navier-Stokes equations are the core of conventional numerical models for atmospheric prediction. These equations are
22 derived from general conservation laws: energy, momentum, mass. Nevertheless, they have an implicit scale invariance
23 symmetry, which is sometimes ignored in regard to other conservation laws (Lovejoy and Schertzer 2013; Palmer 2019). In
24 this work, we exploit this symmetry as the basis for stochastic modelling and prediction of global temperature anomalies.
25 From hourly to centennial time scales, atmospheric fields are characterized by three scaling regimes: at high frequencies the
26 weather, with anomaly fluctuations increasing with the time scale; there is a transition at $\tau_w \sim 10$ days to the macroweather,
27 with fluctuations decreasing with scale; and at low frequencies the climate, again with increasing fluctuations. In recent times,
28 the anthropogenic warming induces the transition to the climate regime at $\tau_c \sim 15$ -20 years, but pre-industrial records show
29 $\tau_c > 100$ years (the Holocene transition scale is still not well known) (Lovejoy 2014). The transition time, τ_w , is the lifetime

30 of planetary structures (Lovejoy and Schertzer 1986, 2010) and is therefore close to the deterministic predictability limit of
31 conventional numerical weather prediction models. This predictability threshold for the models following a deterministic
32 approach is imposed by the high complexity of the system and the sensitive dependence on initial conditions.

33 To extend the predictions to weekly, monthly and seasonal averages, stochasticity is incorporated at different levels to the
34 deterministic prediction systems. The ensemble approach, in which many different “random” realizations are obtained by
35 integrating the model equations from slightly different initial conditions, is fundamentally stochastic. Sampling the attractor
36 of the dynamic system is equivalent to sampling the probability distribution of the possible outputs. Besides this implicit
37 randomness product of chaos, explicit stochastic parameterization schemes are increasingly being incorporated in the
38 prediction systems. Hybrid deterministic-stochastic approaches seem to be the future of macroweather forecasting (Williams
39 2012; Christensen et al. 2017; Davini et al. 2017; Rackow and Juricke 2020). The importance and the current state of stochastic
40 climate modelling has been extensively discussed in the reviews: (Franzke et al. 2015; Palmer 2019).

41 In addition to these stochastic improvements to the deterministic core of conventional Global Circulation Models (GCMs),
42 purely stochastic models have evolved as a complementary approach since the pioneering works of (Hasselmann 1976). For
43 these Hasselmann-type models, the high frequency “weather” is treated as a driving noise of the low frequency components
44 described by integer-order linear ordinary differential equations. The most well-known are the linear inverse models (LIM)
45 (Penland and Matrosova 1994; Penland and Sardeshmukh 1995; Winkler et al. 2001; Newman et al. 2003; Sardeshmukh and
46 Sura 2009). These have been presented as a benchmark for decadal surface temperature forecasts. On the other hand, one of
47 the main limitations of the LIM, is that it implicitly assumes short range exponential temporal decorrelations, while it has been
48 shown that the true decorrelations are closer to long-range power laws (Koscielny-Bunde et al. 1998; Franzke 2012; Rypdal
49 et al. 2013; Yuan et al. 2015). Consequently, LIM models underestimate the memory of the system, imposing a useful limit to
50 the forecast horizon of roughly one year (Newman 2013).

51 In (Lovejoy et al. 2015), the ScaLIng Macroweather Model (SLIMM) was introduced as an alternative stochastic model that
52 respects the scaling symmetry. SLIMM generalizes LIM to use fractional differential equations that involve strong, long-range
53 memories; it is these long-range memories that are exploited in SLIMM forecasts. The solution to the fractional differential
54 equation in SLIMM is a fractional Gaussian noise process that is used to model the natural temperature variability.

55 In a recent series of papers (Lovejoy 2019a, 2020, 2021a, b; Lovejoy et al. 2021), the classical Energy Balance Equation (EBE)
56 is generalized to fractional orders: the Fractional EBE (FEBE). The phenomenological derivation of the FEBE complements
57 derivations based on the classical continuum mechanics heat equation and of the more general Fractional Heat Equation (FHE)
58 (Lovejoy 2020), which is a fractional diffusion equation that has been studied in the statistical physics literature. When the
59 FEBE is driven by a Gaussian white noise, the result is fractional Relaxation noise (fRn) that generalizes the classical Ornstein-
60 Uhlenbeck process and its high-frequency limit is a fractional Gaussian noise process (fGn) that generalizes Brownian motion
61 (Lovejoy 2019a). In that sense, the fractional differential equation and the corresponding fGn solution exploited in SLIMM
62 are the high-frequency approximations of the FEBE and its fRn solution, respectively.

63 In (Del Rio Amador and Lovejoy 2019) (hereafter DRAL) the Stochastic Seasonal to Interannual Prediction System (StocSIPS)
64 was introduced and applied to the prediction of globally averaged monthly temperature in the macroweather regime. StocSIPS
65 includes SLIMM as the core model to forecast the natural variability component of the temperature field, but also represents
66 a more general framework for modelling the seasonality and the anthropogenic trend and the possible inclusion of other
67 atmospheric fields at different temporal and spatial resolutions. In this sense, StocSIPS is the general system and SLIMM is
68 the main part of it dedicated to the modelling of the stationary scaling series. StocSIPS also improves the mathematical and
69 numerical techniques used in the original SLIMM.

70 In DRAL, we presented the basic theory behind StocSIPS and applied it to the prediction of globally averaged series showing
71 verification skill scores in both deterministic and probabilistic modes. We also compared hindcasts with Canada’s operational
72 long-range forecast system, the Canadian Seasonal to Interannual Prediction System (CanSIPS) and we showed that StocSIPS
73 is just as accurate for one-month forecasts, but significantly more accurate for longer lead times.

74 In this paper (specifically in Sections 2.2 and 2.3), we verify that the scaling symmetry, which is the basis of StocSIPS, also
75 holds at the regional level for monthly surface temperature, although some modifications must be introduced in the pre-
76 processing of the tropical ocean temperature anomalies. In Sect. 2.4, we describe these particularities together with some
77 theoretical details, although we purposely placed the most technical aspects in Appendix A, so the main body of the article
78 remains more results-based without too many overwhelming technicalities. Although all the equations and details relevant to
79 this paper are given in the main text or in Appendix A, the interested reader could refer to the more detailed theoretical
80 description given in DRAL. The applicability of the model for all the regional series was confirmed through statistically testing
81 in the second part of Sect. 2.4 and by contrasting the theoretically expected skill scores (if the model were perfect) with actual
82 hindcast verification results for the natural temperature variability in Sect. 3.1. Finally, in Sect. 3.2 we apply StocSIPS to
83 obtain monthly and seasonal predictions of the surface temperature and we show some preliminary comparisons with multi-
84 model ensemble (MME) GCM results.

85 For one month lead time, our simple stochastic model shows similar values of the skill scores than the much more complex
86 conventional models, with the advantage that it is much less expensive computationally and it can be easily adapted to direct
87 hyperlocal prediction without need for downscaling. From a forecast point of view, GCMs can be seen as an initial value
88 problem for generating many “stochastic” realizations of the state of the atmosphere, while StocSIPS is effectively a “past
89 value problem” that estimates the most probable future state from long series of past data. The results obtained validate
90 StocSIPS as a good alternative and a complementary approach to conventional numerical models. This complementarity is the
91 basis for combining the two in a hybrid model that would bring the best of both worlds.

93 2.1 Data preprocessing

94 In this study, the reference observational datasets are monthly average surface temperature (T2m) from the National Centers
 95 for Environmental Prediction/National Center for Atmospheric Research (NCEP/NCAR) Reanalysis 1 (Kalnay et al. 1996;
 96 NCEP/NCAR 2020). The data were accessed on January 3, 2020 and it covers the period January 1948 to December 2019
 97 (864 months in total). All data were interpolated to a 2.5° latitude \times 2.5° longitude grid across the globe for a total of $73 \times$
 98 $144 = 10512$ grid points. Our objective is to model and predict this dataset using the Stochastic Seasonal to Interannual
 99 Prediction System (StocSIPS).

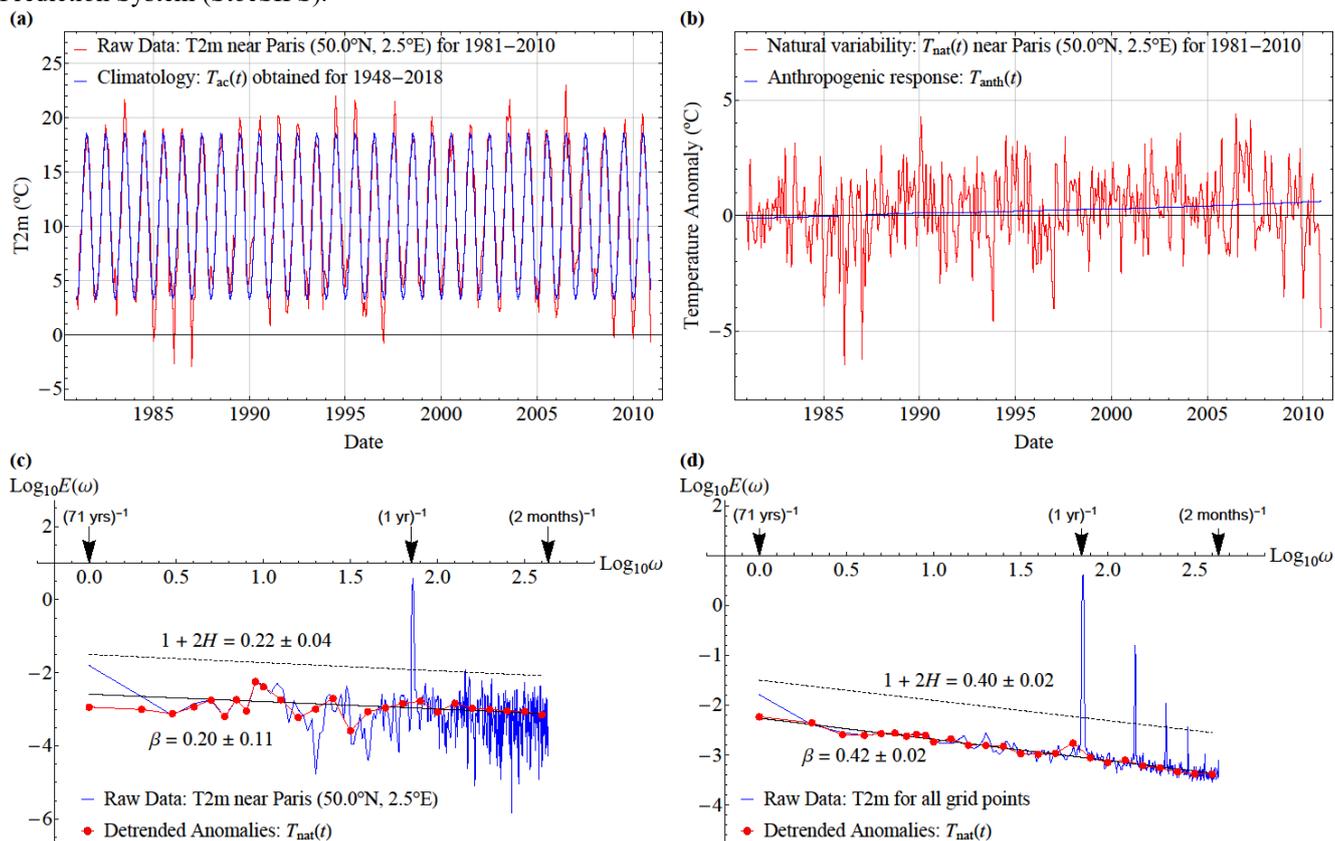


Fig. 1 Example of signal pre-processing and spectra for the grid point with coordinates 50.0°N, 2.5°E (near Paris, France). (a) Raw temperature data, T (in red), and the periodic signal, T_{ac} (in blue). Only the period 1981–2010 is shown for visual clarity. (b) The zero-mean residual natural variability component, T_{nat} and the anthropogenic trend, T_{anth} (red and blue, respectively). (c) Spectra of the raw temperature series and the residual component, T_{nat} (blue and red, respectively). The exponent, β was obtained from the linear regression of the smoothed spectrum. The reference dashed line with slope $1 + 2H$ was also included. (d) Similar to (c), but now considering the average spectra for all the 10512 grid points.

100 StocSIPS was presented in DRAL and applied to the prediction of globally averaged temperature in the macroweather regime.
 101 The main idea behind it is to consider the temperature series at position \mathbf{x} as a combination of three independent signals:

$$102 \quad T(\mathbf{x}, t) = T_{\text{ac}}(\mathbf{x}, t) + T_{\text{anth}}(\mathbf{x}, t) + T_{\text{nat}}(\mathbf{x}, t). \quad (1)$$

103 The first component, $T_{\text{ac}}(\mathbf{x}, t)$, is the periodic annual cycle and is obtained from the mean temperature for each month in some
 104 reference period (here taken as the full length of the temperature datasets: 1948-2019). We assume that, for the time scales
 105 involved in the modelling and prediction problems, the annual cycle is unchanged. Also, for such a long verification period,
 106 the differences with the anomalies obtained using leave-one-out cross-validation methods are negligible. In Fig. 1a, we show
 107 an example of the raw temperature data, T (in red), and the periodic signal, T_{ac} (in blue), for the time series corresponding to
 108 the coordinates 50.0°N, 2.5°E (near Paris, France). In the graph, only the period 1981-2010 is shown for visual clarity.

109 The second component, $T_{\text{anth}}(\mathbf{x}, t)$, is a deterministic low-frequency response to anthropogenic forcings. It can be modelled
 110 as a response to equivalent-CO₂ (CO₂eq) radiative forcing as the one used in CMIP5 simulations (Meinshausen et al. 2011):

$$111 \quad T_{\text{anth}}(\mathbf{x}, t) = \lambda_{2 \times \text{CO}_2 \text{eq}}(\mathbf{x}) \log_2 \left[\rho_{\text{CO}_2 \text{eq}}(t) / \rho_{\text{CO}_2 \text{eq,pre}} \right], \quad (2)$$

112 where $\rho_{\text{CO}_2 \text{eq}}$ is the observed globally-averaged equivalent-CO₂ concentration with preindustrial value $\rho_{\text{CO}_2 \text{eq,pre}} = 277$ ppm
 113 and $\lambda_{2 \times \text{CO}_2 \text{eq}}(\mathbf{x})$ is the transient climate sensitivity at position \mathbf{x} (that excludes delayed responses) related to the doubling of
 114 atmospheric equivalent-CO₂ concentrations. For $\rho_{\text{CO}_2 \text{eq}}$ we used the CMIP5 simulation values (Meinshausen et al. 2011). The
 115 definition of CO₂eq includes not only greenhouse gases, but also aerosols, with their corresponding cooling effect. The
 116 sensitivity $\lambda_{2 \times \text{CO}_2 \text{eq}}(\mathbf{x})$ is estimated from the linear regression of $T(\mathbf{x}, t)$ vs. $\log_2[\rho_{\text{CO}_2 \text{eq}}(t) / \rho_{\text{CO}_2 \text{eq,pre}}]$. This relationship
 117 ignores memory effects, but these are not too strong during periods where the forcing continues to increase. The zero-mean
 118 residual natural variability component, $T_{\text{nat}}(\mathbf{x}, t)$, includes “internal” variability and the response of the system to other natural
 119 forcings (e.g.: volcanic and solar). Both components, T_{anth} and T_{nat} , are shown in Fig. 1b (blue and red, respectively) for the
 120 same point as in Fig. 1a with coordinates 50.0°N, 2.5°E. At this location, it could be argued that the anthropogenic trend is
 121 insignificant compared to the amplitude of the natural component, but at some other locations it is more relevant. Besides, the
 122 cumulative effect of T_{anth} for all the grid points is highly relevant for the globally averaged temperature (see Fig. 5 in DRAL).
 123 Instead of using CO₂eq, alternatively, we could have used the CO₂ concentration in Eq. (2) as a surrogate for all anthropogenic
 124 effects, avoiding various uncertain radiative assumptions needed to estimate CO₂eq (especially aerosols). Nevertheless, from
 125 the point of view of detrending, the residuals, T_{nat} , would remain almost unchanged because of the nearly linear relation
 126 between the actual CO₂ concentration and the estimated equivalent concentration (correlation coefficient > 0.993). There are
 127 more rigorous methods of detrending the original signal to obtain independent components with “stationary” residuals while
 128 preserving the length of the time series [e.g.: empirical mode decomposition (EMD) (Zeiler et al. 2010), ensemble empirical
 129 mode decomposition (EEMD) (Wu and Huang 2009), LOESS (Cleveland and Devlin 1988; Clarke and Richardson 2020)].
 130 Nevertheless, the method used here gives a direct physical meaning to the residual, T_{nat} , and to the low-frequency trend, T_{anth} .
 131 It is also accurate enough for obtaining the detrended temperature anomalies, whose characterization, modelling and prediction

132 are the focus of the following sections. A more accurate method that takes into account the physics of the system adding
 133 memory effects to the heat balance equation, was presented in (Procyk et al. 2020).

134 2.2 Spectra

135 The effects of the detrending in the frequency domain can be observed by comparing the spectra of the raw temperature series
 136 and the residual component, T_{nat} . In Fig. 1c we show these two spectra in a log-log scale in blue and red, respectively, for the
 137 grid point with coordinates 50.0°N, 2.5°E. The spectrum of the detrended series was smoothed by taking averages with
 138 logarithmically spaced bins. Notice that the peak corresponding to the annual cycle was removed along with the signal T_{ac} , as
 139 well as the low-frequency response corresponding to T_{anth} . The frequency, ω , is given in units of 72^{-1} yrs^{-1} (72 years is the
 140 length of the series).

141 After removing the peaks corresponding to the annual cycle (and harmonics) and the low-frequency response, the only relevant
 142 feature of the spectrum of the detrended anomalies, $E(\omega)$, is its scale invariance (power-law behaviour):

$$143 \quad E(\omega) \propto \omega^{-\beta}. \quad (3)$$

144 The exponent, $\beta = 0.20 \pm 0.11$, can be obtained from the linear regression of the smoothed spectrum (shown in red). The line
 145 corresponding to the best fit is shown in black in the figure. We also included a reference dashed line with slope $1 + 2H$,
 146 where H is the fluctuation exponent (see next section).

147 The scaling is even more noticeable in the less noisy spectrum shown in Fig. 1d, obtained by averaging the spectra of all the
 148 10512 grid points. Now the peaks corresponding to the periodic signal and the low-frequency contribution associated with
 149 anthropogenic effects are more clearly visible. The value of the exponent obtained in this case is $\beta = 0.42 \pm 0.02$. The
 150 implications of this scale-invariance will be treated in more detail in the following sections.

151 2.3 Scaling

152 In DRAL, it was shown that, for the case of globally averaged monthly atmospheric surface temperature, the statistics of
 153 $T_{\text{nat}}(t)$ are characterized by one main symmetry: the power-law (scaling) behaviour of the average of the fluctuations, ΔT , as
 154 a function of the time scale, Δt :

$$155 \quad \langle |\Delta T(\Delta t)| \rangle \propto \Delta t^H, \quad (4)$$

156 where H is the fluctuation exponent and the brackets $\langle \cdot \rangle$ denote ensemble averaging. For $-1 < H < 0$, Haar fluctuations, not
 157 differences, should be used (Lovejoy and Schertzer 2012a). Many examples of the low intermittency (“spikiness”) of the
 158 temperature fluctuations are given in (Lovejoy and Schertzer 2013). Equivalently to Eq. (4), in the frequency domain the
 159 spectrum satisfies the previously mentioned equation: $E(\omega) \propto \omega^{-\beta}$, with $\beta = 1 + 2H$ for monofractal processes. These
 160 statistical symmetries are not exclusive to the globally averaged temperature. There are many empirical results that show a
 161 “colored noise” scaling behaviour in local temperature spectra as well as in many other atmospheric variables (Brockwell and
 162 Davis 1991; Blender et al. 2006; Box et al. 2008; Lovejoy and Schertzer 2013; Varotsos et al. 2013; Christensen et al. 2015).

163 For globally averaged temperature at scales between one month and several decades, there is a single scaling regime with $H <$
 164 0. If we analyze temperature time series from daily (or lower) time scales, we find that, in general, there is a transition between
 165 two scaling regimes: from the weather, characterized by fluctuations increasing with the time scale ($H > 0$), to the
 166 macroweather regime where fluctuations tend to cancel out as the time scale increases ($H < 0$).

167 This transition in the statistical properties of the atmosphere at scales of the order of $\tau_w \approx 10$ days, has been theorized by
 168 Lovejoy and Schertzer (1986) as the lifetime of planetary sized structures and estimated from first principles from knowledge
 169 of the solar output and the efficiency of conversion from solar to mechanical energy (Lovejoy and Schertzer 2010). A similar
 170 transition at $\tau_w \approx 1$ year was observed for the average surface temperature over the ocean (Lovejoy and Schertzer 2013).

171 The fluctuation exponents that characterize the weather and the macroweather regimes for air surface temperature (H_w and
 172 H_{mw} , respectively), as well as the transition scale τ_w , are functions of position with a strong dependence on the latitude. In
 173 Fig. 2, we show a map of the exponents obtained from the Haar fluctuation analysis (Lovejoy and Schertzer 2012a) in the
 174 high-frequency scaling regime between 2 months and 2 years. In general, there is a consistent difference between the
 175 macroweather exponents of surface temperature over the oceans and over land with $-1/2 < H_{mw}^{land} < H_{mw}^{ocean} < 0$ (the ocean
 176 is more persistent and the fluctuations cancel out more slowly). Also, for any position over land and for most of the ocean, we
 177 find that $\tau_w < 1$ month, so for surface temperature at monthly resolution, only the macroweather regime is observed. Only for
 178 the tropical ocean we do find a well-defined transition with τ_w as much as 2 years. Consequently, for this region, at time scales
 179 $\Delta t < \tau_w$ the statistics of the fluctuations are those of the weather regime with positive exponents (red in Fig. 2). This longer
 180 transition in the SST corresponds to an analogous “ocean weather” – “ocean macroweather” transition (Lovejoy and Schertzer

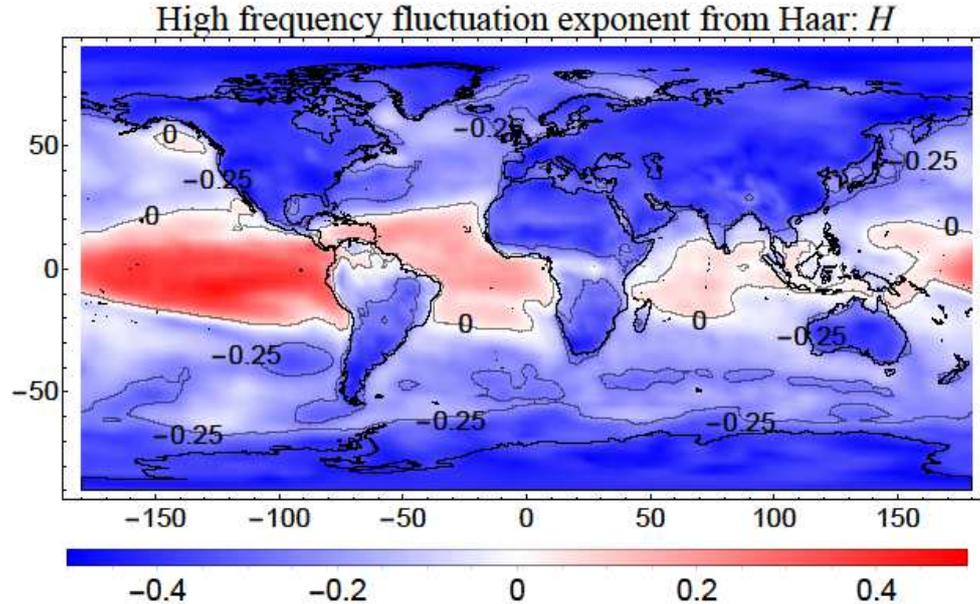


Fig. 2 Map of the fluctuation exponents obtained from the Haar fluctuation analysis (Lovejoy and Schertzer 2012a), in the high-frequency scaling regime between 2 months and 2 years.

181 2012b). It corresponds to lifetimes of large-scale ocean gyres (and other structures) that live much longer than atmospheric
 182 structures.

183 As an example, we show the Haar fluctuation analysis of the time series presented in Fig. 3a. We choose a point over land
 184 (time series in blue in Fig. 3a) with coordinates 50.0°N, 2.5°E (same grid point used before in Sect. 2.1) and a point in the
 185 tropical ocean (red in Fig. 3a) with coordinates 7.5°S, 30°W. In Fig. 3b, we show the average fluctuation as a function of the

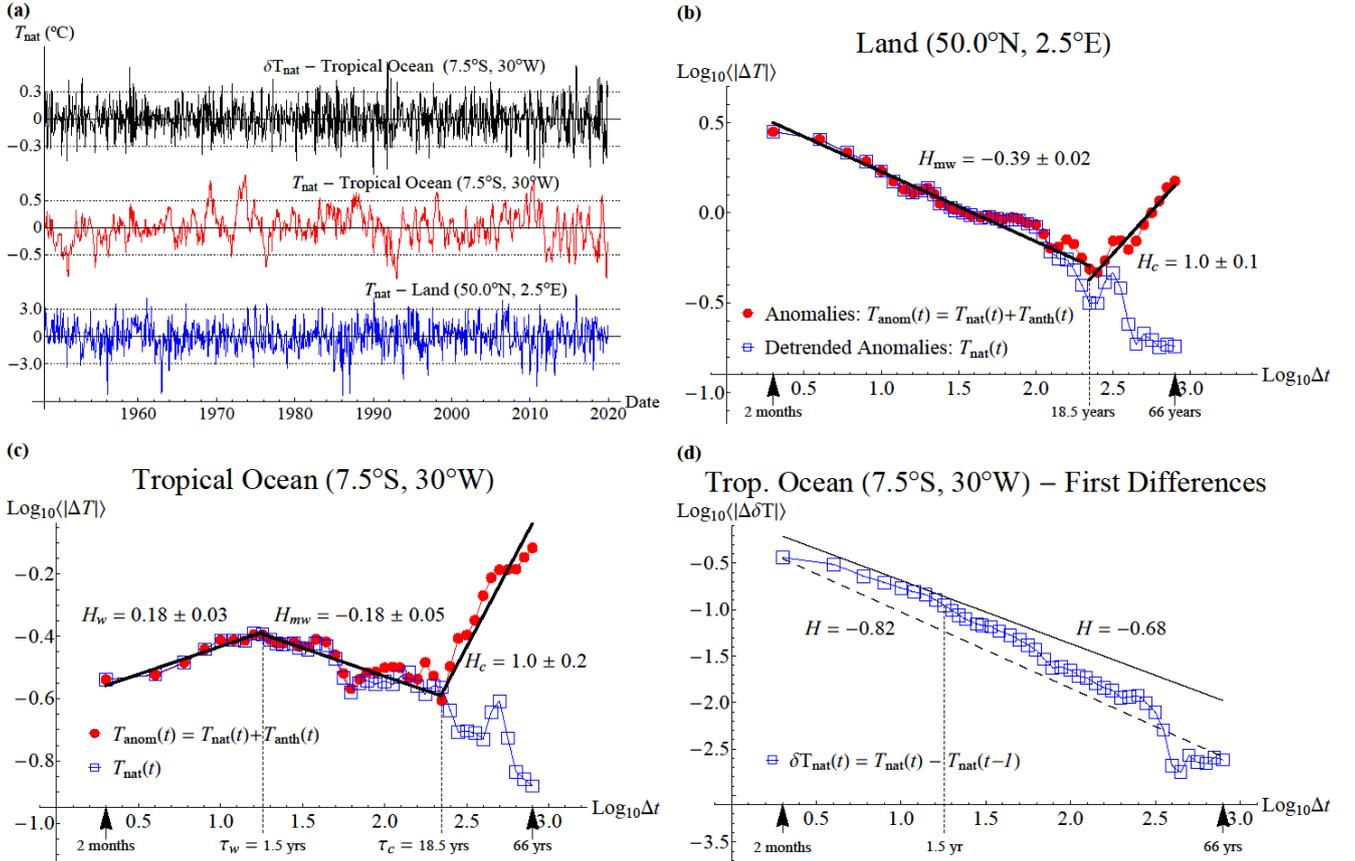


Fig. 3 Examples of Haar fluctuation analysis for two points, one over land and one over ocean. (a) In blue, time series for a point over land with coordinates 50.0°N, 2.5°E (same grid point used before in Sect. 2.1); in red, for a point over ocean located at 7.5°S, 30°W and in black, the series of the temperature differences, $\delta T_{\text{nat}}(t) = T_{\text{nat}}(t) - T_{\text{nat}}(t-1)$, for the same point over ocean (increments of the time series in red). (b) Average fluctuation as a function of the time scale before and after removing the anthropogenic trend for the point over land (red line with circles for the anomalies before removing the anthropogenic component and blue line with empty squares for the detrended anomalies). The reference lines with slopes $H_{mw} = -0.39 \pm 0.02$ and $H_{mw} = 1.0 \pm 0.1$ were obtained from regression of the anomalies' fluctuations in the respective macroweather and climate regimes. (c) Same as in (b) but now for the point over ocean. The three regimes (weather, macroweather and climate) are observed for this point. The corresponding transition scales and the respective exponents obtained from linear regression are also included in the graph. (d) Haar fluctuation analysis of the series of increments $\delta T_{\text{nat}}(t)$ for the point over ocean. The dashed line included as reference has slope $H = H_w - 1 = -0.82$, where H_w is the one shown in (c) and the solid line has a slope $H = -0.68$, which is the exponent obtained from the maximum likelihood method assuming that δT_{nat} is a fractional Gaussian noise (fGn) process (see next section).

186 time scale before and after removing the anthropogenic trend for the point over land [red line with circles for the anomalies
187 before removing the anthropogenic component ($T_{\text{anom}} = T_{\text{nat}} + T_{\text{anth}}$) and blue line with empty squares for the detrended
188 anomalies (T_{nat})]. The reference line with slope $H_{mw} = -0.39 \pm 0.02$ was obtained from regression of the residuals'
189 fluctuations between 2 months and 18.5 years. The units for Δt and ΔT are months and °C, respectively.

190 Notice that the anthropogenic warming breaks the scaling of the undetrended anomalies' fluctuations at a time scale of 15-20
191 years (the fluctuations start to increase with the scale at ~ 200 months). The fluctuation exponent for this low-frequency
192 (climate) regime is $H_c = 1.0 \pm 0.1$ – i.e., the fluctuations increase linearly with time following the almost linear growth of
193 CO₂ concentration in recent epochs. The residual natural variability, on the other hand, shows reasonably good scaling for the
194 whole period analyzed (66 years). In analysis of temperature records from preindustrial multiproxies and GCMs preindustrial
195 control runs (Lovejoy 2014), evidence was presented showing that the range of scaling with decreasing fluctuations (pre-
196 industrial macroweather) may extend to more than 100 years.

197 As we mentioned before, for this point over land, only one regime with fluctuations decreasing with the time scale (the
198 macroweather regime) is present for the natural variability. So, we can conclude that, at this location, the weather-
199 macroweather transition occurs at $\tau_w < 1$ month (maximum resolution of the analyzed data). This was confirmed using 6-
200 hours resolution data. In contrast, as we show in Fig. 3c, if we analyze the grid point in the tropical region over the ocean (time
201 series in red in Fig. 3a), there is a clear transition at $\tau_w \sim 1.5$ years from the weather regime (with fluctuations increasing with
202 the scale) to the macroweather regime (with decreasing average fluctuations). A further transition occurs in the undetrended
203 anomalies at $\tau_c \sim 18.5$ years to the climate regime, where fluctuations start to increase again with the time scale. As before,
204 this transition in recent epochs is induced by anthropogenic effects. The actual transition in the natural variability, as obtained
205 from preindustrial temperature records, apparently occurs at time scales longer than 100 years, which is consistent with the
206 blue curves for T_{nat} in Figs. 3b and 3c after we remove the anthropogenic trend. The fluctuation exponent for the three regimes,
207 weather-macroweather-climate, has values $H_w = 0.18 \pm 0.03$, $H_{mw} = -0.18 \pm 0.05$ and $H_c = 1.0 \pm 0.2$, respectively
208 (shown in the graph) consistent with a smooth low-frequency behaviour.

209 A visual comparison between the blue and red curves in Fig. 3a shows a clear difference in the temperature behaviour at these
210 two grid points. While over land, consecutive values of temperature tend to cancel out, over the ocean the temperature is more
211 persistent and only after several time steps the anomalies change sign. This is confirmed in the Haar fluctuation analysis shown
212 in Figs. 3b and 3c. This difference in the statistical behaviour imply that, while a fractional Gaussian noise (fGn) model is a
213 good fit for the extratropics, we cannot use it to describe the tropical region. Nevertheless, if we take the first differences in
214 the time series for the grid point over the tropical ocean, the new series $\delta T_{\text{nat}}(t) = T_{\text{nat}}(t) - T_{\text{nat}}(t - 1)$ (shown in black in
215 Fig. 3a) has a statistical behaviour which is clearly more similar to the series over land with consecutive fluctuations cancelling
216 out. As we can see in the graph shown in Fig. 3d, the new series $\delta T_{\text{nat}}(t)$ has a scaling regime for small Δt with negative
217 fluctuation exponent similar to that of Fig. 3b. By taking first differences in the tropics, we are able to use fGn process
218 everywhere to predict the time series, then we can go back to the original series for those places by taking cumulative sums.

219 There is still a change in the slope at $\tau_w \sim 1.5$ years, corresponding to the one in the original series shown in Fig. 3c. The dashed
220 line included as reference has a slope $H = H_w - 1 = -0.82$. The series δT_{nat} , being the increments of the series T_{nat} , should
221 have an exponent of the dominant high frequencies reduced by one. We also included in solid black, a reference line with slope
222 $H = -0.68$, which is the exponent obtained from the maximum likelihood method assuming that δT_{nat} is an fGn process (see
223 Sect. 2.4.2).

224 These examples – shown here for two different positions – are representative of the behaviour of the natural temperature
225 variability all over the Earth. In fact, by taking the first differences of the time series in those places over the tropical ocean
226 with weather regime at monthly resolution, we can reduce our analysis to only one case of self-similar time series with negative
227 exponent in the range $-1 < H < 0$. This simplification emphasizes the role of the scaling symmetry, which is sometimes
228 ignored in regard to other conservation laws, in spite of being also present in the Navier-Stokes equations (Lovejoy and
229 Schertzer 2013; Palmer 2019), which are the core of conventional numerical models for atmospheric prediction and hence
230 respected by them. In this work, we exploit this symmetry as the basis for stochastic modelling and prediction of global
231 temperature anomalies.

232 2.4 Stochastic modelling using fGn and fRn

233 2.4.1 Properties of fGn, fRn

234 Together with the scaling symmetry presented in the previous section, we also assume the Gaussianity of the natural
235 temperature variability. This Gaussian hypothesis was verified in DRAL for globally averaged monthly temperature in the
236 macroweather regime. Although the Gaussian assumption is commonly made, it is worth underlining that it is somewhat
237 surprising that it is a reasonable model for macroweather time series. Recall that Gaussian statistics imply that macroweather
238 in time has little or no intermittency (the series are mono-, not multifractal, the transitions are not “spiky”). This contrasts
239 with macroweather in space, which is highly intermittent, as well as the existence of highly intermittent, nonGaussian,
240 multifractal spatial and temporal statistics in the weather and climate regimes (Lovejoy 2018).

241 The scaling of the temperature fluctuations and spectrum implies that there are power-law correlations in the system and hence
242 a large memory effect that can be exploited. In (Lovejoy 2019a; Lovejoy et al. 2021), it was argued that the origin of this
243 memory are the Earth’s hierarchical, scaling energy storage mechanisms whereby anomalies in energy fluxes either external
244 (e.g. anthropogenic) or internal can be stored for long periods. It was argued that to a good approximation, the temperature
245 satisfies the Fractional Energy Balance Equation (FEBE) that has a high-frequency scaling storage term and a low-frequency
246 thermal equilibrium term. When the FEBE is internally forced by a Gaussian white noise, the temperature response is the
247 statistically stationary fractional Relaxation noise (fRn) process (Lovejoy 2019b).

248 However, at time scales shorter than the relaxation time (of the order of a few years), the (scaling) storage term is dominant
249 and, for exponents $-1/2 < H < 0$, the temperature response is a fractional Gaussian noise (fGn) process. This was the

250 approximation made in DRAL and is empirically valid for all land areas and most of the oceans. The exceptions are some parts
251 of the tropical ocean where $0 < H < 1$ (Figs. 2 and 4a), we return to these below.

252 The original idea of modelling the natural variability using an fGn process was presented in (Lovejoy et al. 2015) as the
253 ScaLIng Macroweather Model (SLIMM). In DRAL, StocSIPS was introduced as a general system that includes SLIMM as
254 the core prediction model. StocSIPS also improves the mathematical and numerical techniques of SLIMM. It was applied to
255 the prediction of globally averaged temperature series since 1880. The comparison of StocSIPS hindcasts with Canada's
256 operational long-range forecast system, the Canadian Seasonal to Interannual Prediction System (CanSIPS), showed that
257 StocSIPS is just as accurate for one-month forecasts, but significantly more accurate for longer lead times.

258 In this paper we extend the globally averaged version of StocSIPS for the prediction of a single temperature time series to the
259 prediction of the full space-time temperature field. The basic theory for fGn processes, used here to model those places where
260 $-1/2 < H < 0$ (most of the planet), is summarized in Appendix A. An fGn process is fully characterized by two parameters
261 (assuming zero mean): the fluctuation exponent, H , and the standard deviation, σ_T .

262 We mentioned that for some tropical ocean regions, $0 < H < 1$. While these may still be modelled by fRn processes, the high-
263 frequency approximation to fRn is no longer an fGn process, but rather a fractional Brownian motion (fBm) process, and we
264 must use the correlation function for fRn given in Appendix Aiii., Eq. (A9). For those regions with positive H , the first
265 differences of the temperature, $\delta T_{\text{nat}}(t) = T_{\text{nat}}(t) - T_{\text{nat}}(t - 1)$, has H values reduced by 1, so for δT_{nat} we also have $-1 <$
266 $H < 0$. That is, either the natural temperature variability itself or its first differences can be modelled by an fGn process. In
267 those places where $H > 0$ for the high frequencies, it would be equivalent to modelling them with an fBm or fRn process. Of
268 course, a true fBm would only have one scaling regime with positive fluctuation exponent, instead of the bi-scaling regime
269 shown for the detrended anomalies in Fig. 3c. To model those series as an fBm process is an approximation that would work
270 well for the high frequencies, but that would fail in reproducing the low frequency behaviour.

271 **2.4.2 Parameter estimates and model adequacy**

272 With the distinction in the tropical region where we take the first differences to adjust everything to an fGn model, we conclude
273 that, to model the actual temperature field for the globe (including the anthropogenic trend), for each grid point of the
274 NCEP/NCAR Reanalysis 1 data we may estimate the three parameters (H , σ_T , $\lambda_{2 \times \text{CO}_2 \text{eq}}$). For the first two, we use the
275 maximum likelihood method described in Appendix 1 of DRAL and for the sensitivity we use the regression described in Sect.
276 2.1. To verify the model adequacy, we use Eq. (A8) to obtain the residual innovations, $\gamma(t)$, then we obtain its variance, σ_γ ,
277 and its fluctuation exponent, H_γ , using the maximum likelihood method; they should be equal to 1 and $-1/2$, respectively
278 (white noise processes are particular cases of fGn with $H = -1/2$). The results are summarized in Fig. 4.

279 A map of the maximum likelihood estimates of the temperature fluctuation exponent is shown in Fig. 4a. These values are
280 more accurate and give a better fit of our model than the high-frequency Haar estimates shown in Fig. 2. Notice that for most
281 of the globe and all of the land, the values are in the range $-1/2 < H < 0$, which is characteristic of long-range memory fGn

282 processes with nonsummable correlation functions, i.e. the sum over Δt of the series with elements given by Eq. (A3) diverges
 283 for this range of H . There is a discontinuity from negative to positive values of H as we approach the tropical ocean,
 284 corresponding to the change in model from fGn to fBm (or, equivalently, from the description as an fGn of the natural
 285 temperature variability, T_{nat} , to the description of the temperature differences, δT_{nat}). In most of the tropical ocean (red regions
 286 in the map), the natural temperature variability has fluctuation exponents in the range $0 < H < 1/2$, whose fBm
 287 approximation has “anti-persistent” increments (consecutive increments are negatively correlated). Only in the eastern
 288 equatorial Pacific (yellow region in the map), do we obtain fluctuation exponents in the range $1/2 < H < 1$, whose fBm
 289 approximation has persistent (positively correlated) increments. It is significant that it is precisely this more predictable region

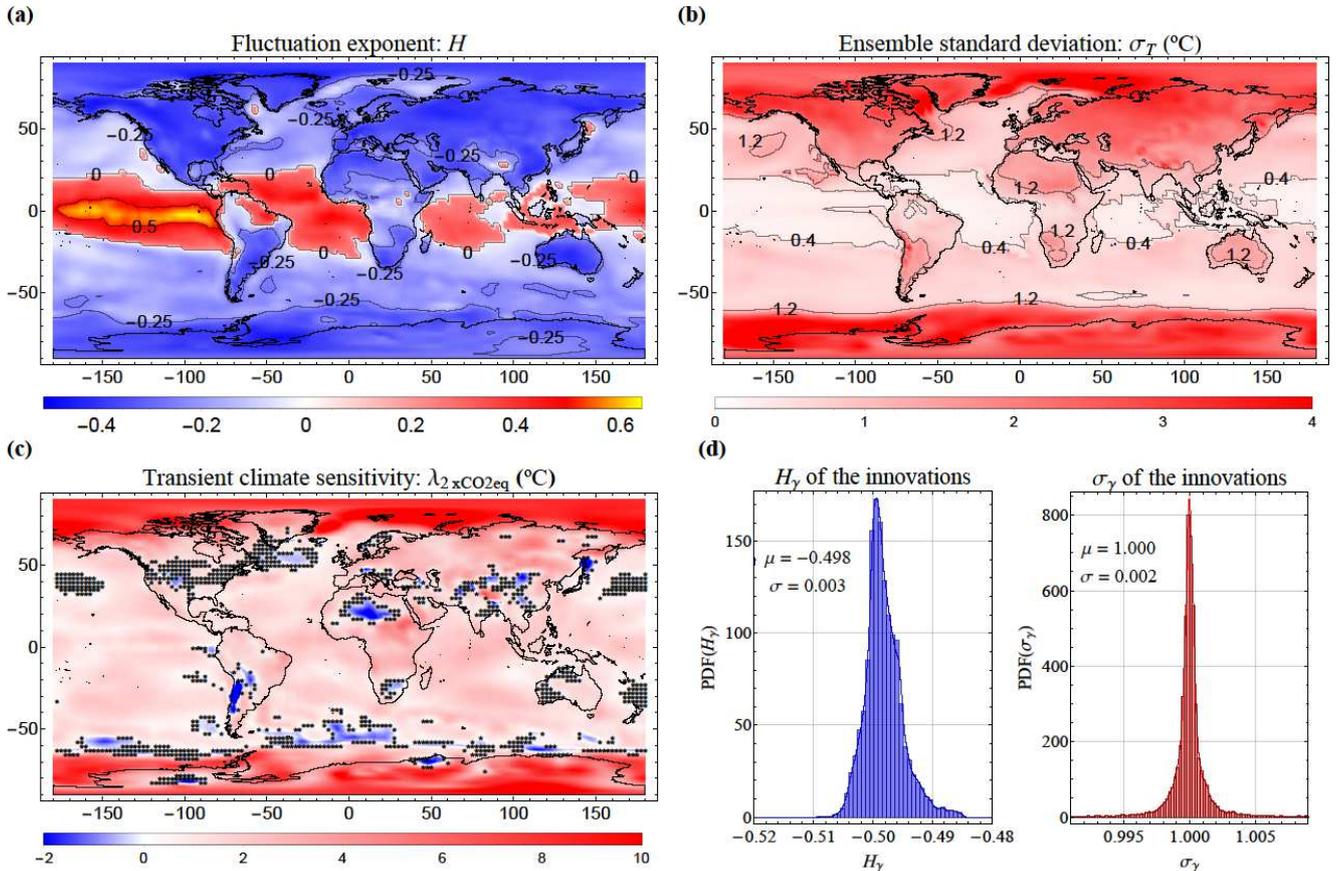


Fig. 4 Estimates of the three parameters (H , σ_T , $\lambda_{2 \times \text{CO}_2 \text{eq}}$) obtained for each grid point and statistics of the innovations, $\gamma(t)$. (a) Maximum likelihood estimates of the temperature fluctuation exponent (compare with the estimates shown in Fig. 2). There is a discontinuity from negative to positive values of H as we approach the tropical ocean, corresponding to the change in model from fGn to fBm. (b) The standard deviation, σ_T , of the infinite ensemble fGn process. (c) Map of the transient climate sensitivity, defined in Eq. (2). The places marked with “*” indicate pixels where the null hypothesis, $\lambda_{2 \times \text{CO}_2 \text{eq}} = 0$, cannot be rejected with more than 90% confidence. (d) Histograms of the fluctuation exponent and the standard deviation of the innovations (H_γ and σ_γ , respectively) for the 10512 grid points. From the histograms, we can conclude that the innovations are very close to white noise for the whole planet ($H_\gamma = -0.498 \pm 0.003$ and $\sigma_\gamma = 1.000 \pm 0.002$).

290 that is associated with the ENSO phenomenon (Trenberth 1997), the strongest interannual signal of climate variability on
 291 Earth.

292 In Fig. 4b we show the values of the parameter σ_T . Although this is the standard deviation of the infinite ensemble fGn process,
 293 for a given finite realization it does not coincide with the usual estimate based on the temporal average:

$$294 \quad SD_T^2 = \frac{1}{N} \sum_{t=1}^N [T_{\text{nat}}(t) - \bar{T}_N]^2, \quad (5)$$

295 where $\bar{T}_N = \sum_{t=1}^N T_{\text{nat}}(t)/N$ (the over-bar notation is used in to denote averaging in time). The biased estimate SD_T ignores
 296 correlations, that are however considered in the maximum likelihood estimate of σ_T . The relation between the two values for
 297 fGn processes depends on the length of the time series and the fluctuation exponent, H , and is given by:

$$298 \quad SD_T^2 = \sigma_T^2 (1 - N^{-2H}) \quad (6)$$

299 (see Sect. 3.3 and Appendix 1 of DRAL). Notice that there is also a discontinuity in the map of σ_T for the same reasons
 300 explained previously. In general, the amplitude of the fluctuations is larger over land than over the ocean; the surface
 301 temperature over the ocean is less variable as this has a higher thermal inertia than land.

302 A map of the transient climate sensitivity, defined in Eq. (2), is shown in Fig. 4c. The places marked with “*” indicate grid
 303 boxes where the null hypothesis, $\lambda_{2 \times \text{CO}_2 \text{eq}} = 0$, cannot be rejected with more than 90% confidence. Notice that these values
 304 depend on the reference dataset. In our case we used the NCEP/NCAR Reanalysis 1, which only has data since 1948. More
 305 precise estimates of the climate sensitivity were obtained by Hébert and Lovejoy (2018) using five observational datasets since
 306 1880. In this paper, we are not aiming at an accurate study of the climate sensitivity. We should consider the values of $\lambda_{2 \times \text{CO}_2 \text{eq}}$
 307 reported here as a parameter used for detrending the temperature time series related to the anthropogenic effects.

308 Finally, in Fig. 4d, we show histograms of the fluctuation exponent and the standard deviation of the innovations (H_γ and σ_γ ,
 309 respectively) for the 10512 grid points. From the histograms, we conclude that the innovations are very close to white noise
 310 for all the places in the planet ($H_\gamma = -0.498 \pm 0.003$ and $\sigma_\gamma = 1.000 \pm 0.002$). So, with a high degree of accuracy, all the
 311 innovation series can be considered NID(0,1), and we can conclude that the fGn model is a good fit to the natural temperature
 312 variability (or its increments in the red and yellow places of the map in Fig. 4a).

313 **3 Results**

314 **3.1 Natural variability forecast**

315 **3.1.1 Model validation through hindcast**

316 In the previous section, we validated the fGn model as a good fit to the natural temperature variability (or to its increments)
 317 by checking the whiteness of the residual innovations. The goal of this section is to further validate the model by using the
 318 theory presented in Appendix Aiv to hindcast only the natural variability – not the anthropogenic signal or the annual cycle –

319 and seeing how well it performs. We test the assumptions made in the model by comparing the theoretically expected skill
 320 scores (expected values if the model were perfect) with the actual scores obtained from hindcasts. All the verification metrics
 321 used in this paper are detailed in Appendix B.

322 Series of hindcasts at monthly resolution, were produced for forecast horizon from 1 to 12 months, in the period of verification
 323 (POV) from December 1950 to November 2019 (the verification starts in December in order to have the same number of
 324 conventional seasons: DJF, MAM, JJA, SON). In this 69-year verification period, each month was independently predicted
 325 using the information available m months before. For each horizon, k , we used a memory $m = 20$ months. For example, to
 326 predict the average temperature for December 1950 with $k = 1$ month, we used the previous 21 months, including November
 327 1950, and the same was done for every verification date up to November 2019 and for all horizons up to $k = 12$ months. The
 328 dependence with the horizon of many scores [e.g. the root mean square error (RMSE)], is obtained from the difference between
 329 hindcasts series at a fixed k and the corresponding series of observations.

330 It is important to point out that the predictor $\hat{T}_{\text{nat}}(t + k)$ (see Eq. (A10)) only depends on the previous $m + 1$ months, from
 331 $T_{\text{nat}}(t - m)$ to $T_{\text{nat}}(t)$, weighted by coefficients that only depend on the fluctuation exponent H (see Fig. 4a). The estimates
 332 of H are quite robust and only small variations were obtained for different training periods, as long as the length of the training
 333 periods is larger than one third of the full length of the time series. Also, only small changes on the skill were appreciated for
 334 small variations in H . In that sense, given the stability on the estimates of the fluctuation exponent, we can use almost all the
 335 observational period for verification leaving only a few months before the first initialization date to use as memory. In all
 336 cases, the observational and forecast anomalies used for verification were calculated in the leave-one-out cross validation
 337 mode.

338 Root Mean Square Error (RMSE)

339 The infinite ensemble expectation of the RMSE is given in Appendix Aiv (Eq. (A13)). This analytical expression is only a
 340 function of the model parameters and does not include any observational data. It is the theoretical value on what the RMSE

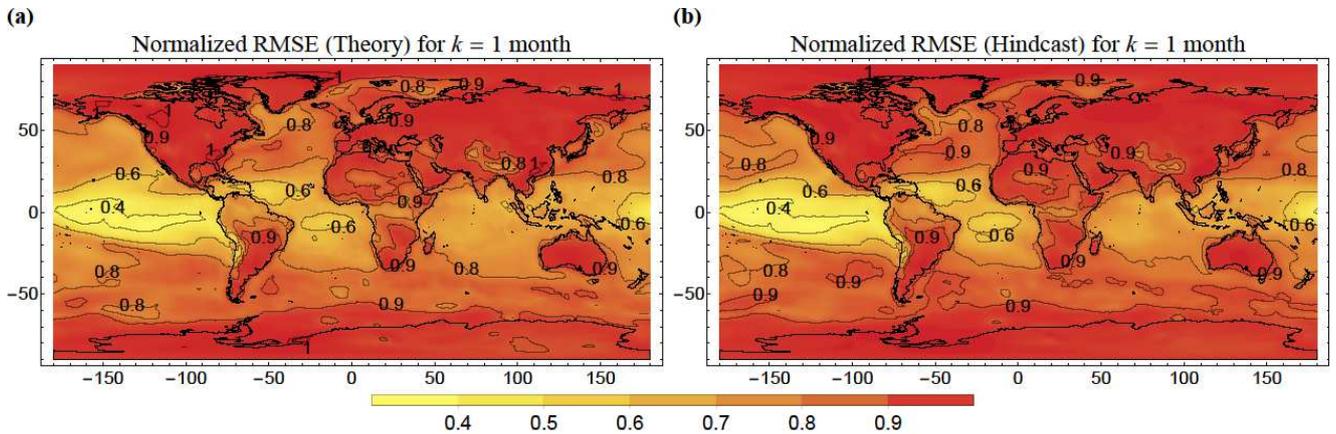


Fig. 5 Theoretical and hindcasts NRMSE for $k = 1$ month. The corresponding RMSEs were obtained using Eqs. (A13) and (B3), respectively, and the normalization standard deviation from Eq. (5) for the natural variability.

341 would be if the model were perfect. To confirm the validity of the theoretical framework for the prediction of the natural
 342 variability component, we compare these expected values for each grid point with the actual verification RMSE obtained from
 343 hindcasts in the POV from December 1950 to November 2019. The all-month verification score for horizon k is obtained
 344 using Eqs. (B2) and (B3) with $N = 828$ months and $T_{\text{nat}}(t + k)$ and $\hat{T}_{\text{nat}}(t + k)$ being the zero mean detrended observational
 345 and predicted anomalies, respectively.

346 The comparison between the theoretical and the actual (obtained from hindcasts) normalized root mean square error (NRMSE)
 347 is shown in Fig. 5 for horizon $k = 1$ month. The NRMSE is the RMSE normalized by the observed standard deviation (Eq. (5)
 348 for the natural variability). The NRMSE may vary from zero to infinity, with lower NRMSE values indicating more skillful
 349 forecasts. NRMSE values greater than 1 indicate that forecasts are less skillful than the climatological average value of the
 350 series. As we pointed out, the theoretical RMSE only depends on the parameters σ_T and H . In general, there is very good
 351 agreement between theory and verification results. The maximum difference between the two maps in Fig. 5 is lower than

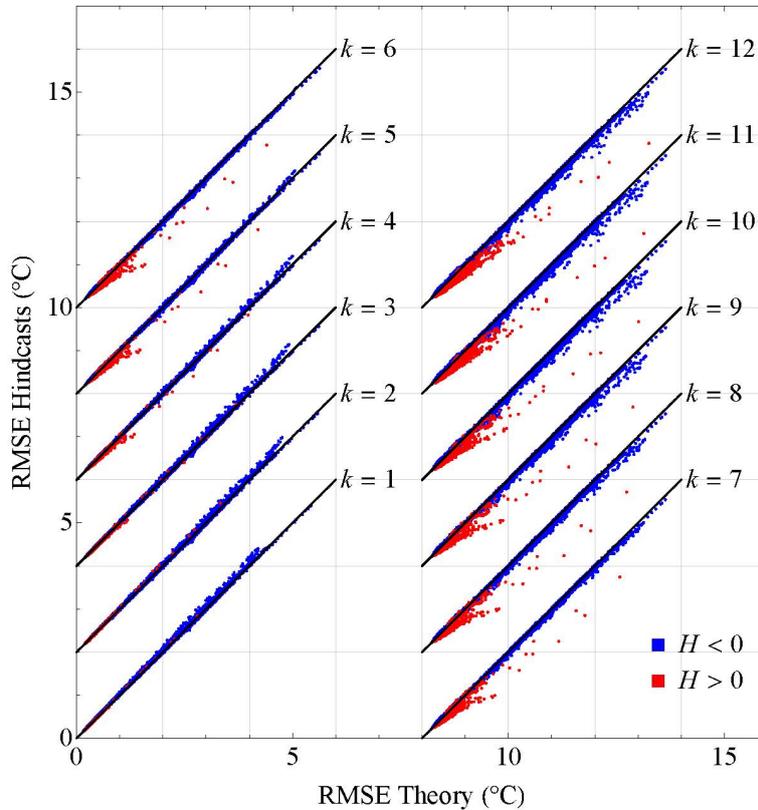


Fig. 6 Scatter plots for each horizon including the 10512 grid points, showing the verification RMSE obtained from hindcasts vs. the expected theoretical $\text{RMSE}_{\text{nat}}^{\text{theory}}$ predicted by Eq. (A13). The graphs were displaced vertically by 2°C (plus a horizontal displacement of 8°C for $k \geq 7$ months) for visual clarity. The black line at 45° is a reference indicating perfect agreement between theory and verification results. The blue points represent locations where $H < 0$ and the natural variability is modeled as an fGn process and the red points are for places where $H > 0$ and we use the fBm model.

352 0.07. The forecast skill is higher over ocean than over land and takes the highest values over the tropical ocean, which
 353 corresponds to the spatial distribution of H values shown in Figs. 2 and 4a.
 354 The maps in Fig. 5 were obtained for $k = 1$ month, but similar maps can be obtained for all forecast horizons from 1 to 12
 355 months. The results of the comparison can be summarized in the scatter plots shown in Fig. 6. The graphs include the 10512
 356 grid points, showing the verification RMSE obtained from hindcasts vs. the expected theoretical $\text{RMSE}_{\text{nat}}^{\text{theory}}$ predicted by Eq.
 357 (A13) for each horizon. As expected, the agreement between the theoretically expected scores and the hindcasts results
 358 decreases as the horizon increases, but it remains quite accurate in all cases with a correlation coefficient larger than 0.998.
 359 For the regions where $H > 0$, the fBm fit is less accurate; however, recall that in those places the actual statistics of the
 360 fluctuations are bi-scaling, while the fBm model assumes a perfectly scaling process. The accuracy of the theory decreases as
 361 the horizon approaches the transition time, τ_w .

362 Mean Square Skill Score (MSSS)

363 Related to the RMSE score, the MSSS is a commonly used metric (see Eq. (B5)). The guidelines of the World Meteorological
 364 Organization (WMO) Standard Verification System for Long-Range Forecasts (LRFs) (WMO 2010a), suggests the MSSS as
 365 a metric for deterministic forecasts (based on the ensemble mean). For leave-one-out cross-validated data in the POV (WMO
 366 2010a), the mean square error (MSE) of the reference climatology forecasts (including the deterministic anthropogenic trend
 367 forecast) is:

$$368 \quad \text{MSE}_c = \left(\frac{N}{N-1} \right)^2 SD_T^2 \quad (7)$$

369 (see Eq. (B4)), where SD_T^2 is the variance of the detrended anomaly series (natural variability component). The MSSS for
 370 horizon k for the natural variability forecast is:

$$371 \quad \text{MSSS}_{\text{nat}}(k) = 1 - \frac{\text{MSE}_{\text{nat}}(k)}{\left(\frac{N}{N-1} \right)^2 SD_T^2}, \quad (8)$$

372 where MSE_{nat} is obtained using Eqs. (B2) with $N = 828$ months and $T_{\text{nat}}(t+k)$ and $\hat{T}_{\text{nat}}(t+k)$ being the zero mean
 373 detrended observational and predicted anomalies, respectively.

374 One consequence of the memory effects in the natural variability is the increase of SD_T^2 with the length of the verification
 375 period given by Eq. (6). This implies that some metrics, such as the MSSS or the NRMSE, will actually have the same
 376 dependence with the duration of the verification period. The longer the verification period, the higher the value of MSSS (lower
 377 for NRMSE), even for the same prediction system. Comparisons between skill scores of different models should always be
 378 made for the same POV (or at least the same length of the POV). As the number of months used for verification increases,
 379 $SD_T^2 \rightarrow \sigma_T^2$ and the MSSS approaches the asymptotic value (determined by H). This effect is small for most values of H , but
 380 is significant if too short verification periods are used or if H is close to zero (e.g.: $SD_T^2/\sigma_T^2 \approx 0.6$ for $N = 100$ months and
 381 $H = -0.1$). See Fig. 9 in DRAL for an example in monthly globally averaged temperature.

382 **Temporal Correlation Coefficient (TCC)**

383 The TCC is another commonly used verification score for deterministic forecasts (see Eq. (B6)). For the natural variability
 384 forecast, the TCC for horizon k is:

385
$$\text{TCC}_{\text{nat}}(k) = \frac{\overline{T_{\text{nat}}(t+k)\hat{T}_{\text{nat}}(t+k)}}{SD_T \sqrt{\overline{\hat{T}_{\text{nat}}(t)^2}}}, \quad (9)$$

386 where the overbars indicate temporal average for a constant k .

387 For the natural variability forecast, the autoregressive coefficients in our predictor were obtained as analytical functions of
 388 only the fluctuation exponent, H (see Eqs. (A10) and (A11)). As we showed in Appendix Biii, if our model is adequate for
 389 describing the natural temperature variability, then the following relationship between the verification TCC_{nat} and MSSS_{nat}
 390 should be satisfied for $k = 1$ month:

391
$$\text{TCC}_{\text{nat}}(1) \approx \sqrt{\text{MSSS}_{\text{nat}}(1)}. \quad (10)$$

392 It does not hold for all horizons in the tropical region due to the use of the fBm rather than fGn model.

393 In Fig. 7 we show maps of the TCC_{nat} and the absolute difference $|\text{TCC}_{\text{nat}} - \sqrt{\text{MSSS}_{\text{nat}}}|$ obtained from hindcasts for $k = 1$
 394 month. The color scale in (b) was rescaled 100 times with respect to (a) so the differences could be perceptible. They are
 395 negligible compared to the values in (a). The maximum differences in Fig. 7b is almost always lower than 0.01 (mean value
 396 of 0.001), which corroborates the adequacy of the fGn model to describe the natural variability.

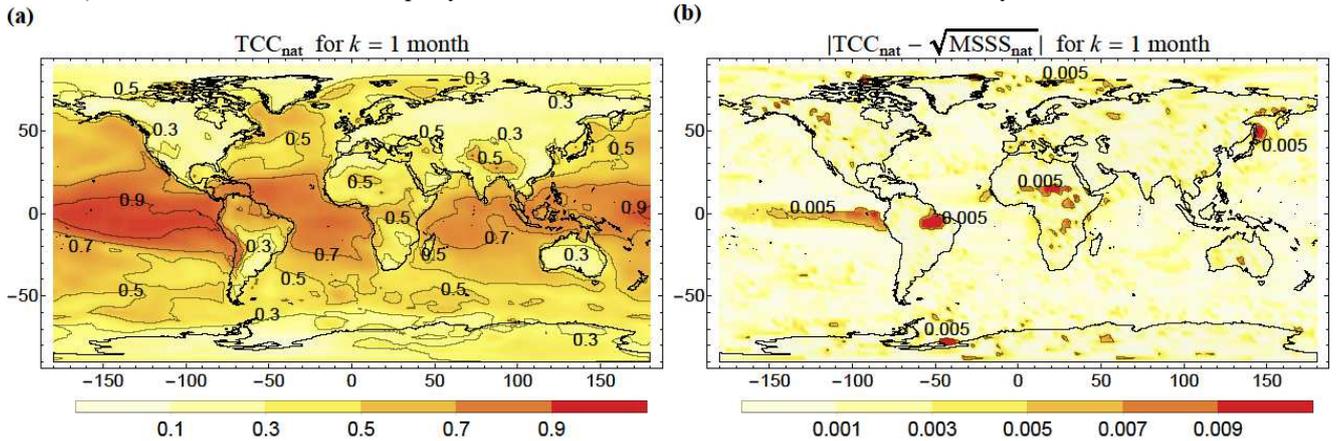


Fig. 7 Maps of TCC_{nat} and the absolute difference $|\text{TCC}_{\text{nat}} - \sqrt{\text{MSSS}_{\text{nat}}}|$ obtained from hindcasts for $k = 1$ month. The colour scale in (b) was rescaled 100 times with respect to (a) so the differences could be perceptible.

397 **3.1.2 Probabilistic scores and reliability**

398 All the skill scores discussed above are recommended by the WMO for assessing deterministic prediction of long-range
 399 forecasts (WMO 2010b). These forecasts are deterministic in the sense that only the ensemble mean is considered, disregarding
 400 the ensemble variance, or more accurately, the prediction of the probability distribution. In this study we only focus on

401 deterministic predictions (deterministic in the previously mentioned sense, recall that we use a stochastic model) because,
 402 given a Gaussian approximation of a probability distribution function, the skill of probabilistic forecasts is mainly dependent
 403 upon the skill of ensemble mean predictions and much less upon predictions of ensemble variances (Kryjov et al. 2006). In
 404 fact, in DRAL it was shown that, assuming a Gaussian distribution for the errors, the Continuous Ranked Probability Score
 405 (CRPS) (Hersbach 2000; Gneiting et al. 2005), which is a commonly used metric for probabilistic forecasts, is related to the
 406 RMSE by:

$$407 \quad \text{CRPS}(k) = \frac{\text{RMSE}(k)}{\sqrt{\pi}} \left[\sqrt{2(1+\text{ESS})} - \sqrt{\text{ESS}} \right], \quad (11)$$

408 where:

$$409 \quad \text{ESS} = \frac{\overline{\sigma_{\text{ensemble}}^2}}{\text{MSE}} \quad (12)$$

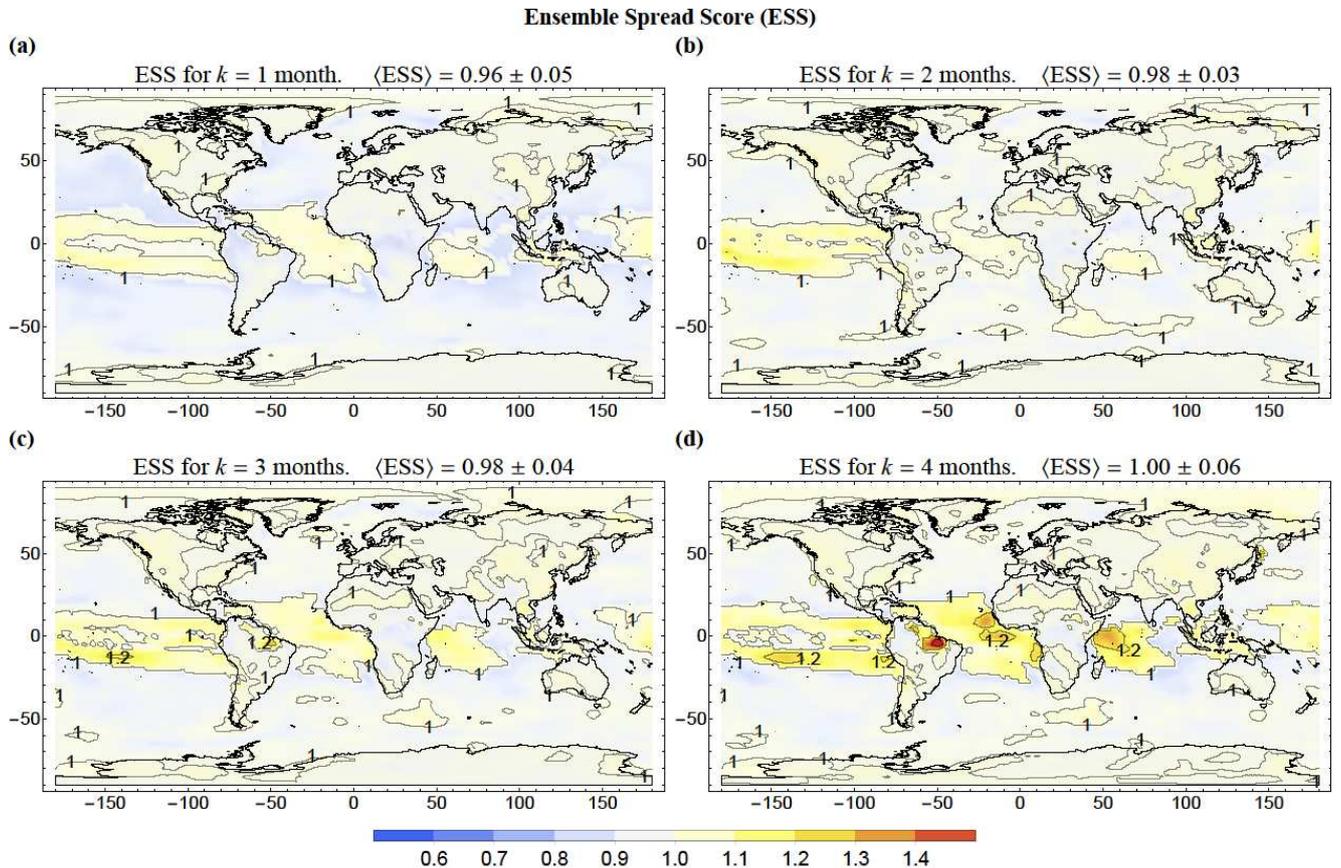


Fig. 8 Maps of ESS of StocSIPS for horizons k from 1 to 4 months (panels (a) to (d), respectively). The values of the ESS are very close to 1, with the exception of the tropical ocean where it tends to be “overdispersive” ($\text{ESS} > 1$). The average values for the globe with one standard deviation are shown in brackets in the map labels.

410 is the ensemble spread score, defined as the ratio between the temporal mean of the intra-ensemble variance, $\overline{\sigma_{\text{ensemble}}^2}$, and
 411 the mean square error between the ensemble mean and the observations (Palmer et al. 2006; Keller and Hense 2011; Pasternack
 412 et al. 2018). The ESS is a commonly used metric to evaluate the reliability of the probabilistic forecast of an ensemble model.
 413 For the case of StocSIPS, which by definition is a Gaussian model with ensemble spread $\sigma_{\text{ensemble}} = \text{RMSE}_{\text{nat}}^{\text{theory}}$ (given by
 414 Eq. (A13)), the agreement between $\text{RMSE}_{\text{nat}}^{\text{theory}}$ and RMSE_{nat} (summarized in Fig. 6 for all horizons) implies that $\text{ESS} \approx 1$
 415 almost everywhere.

416 The graphs shown in Fig. 6 are analogous to spread-error scatterplots (Leutbecher and Palmer 2008). In our case, each point
 417 represents the ensemble spread and the temporal average RMSE for each pixel, instead of the spatially averaged values shown
 418 in Fig 4 of (Leutbecher and Palmer 2008). We could group up and average the values in equally populated bins to produce
 419 more similar spread-error plots, but as they all fall near to the reference diagonal, the conclusions would remain the same.
 420 Other measures used to assess the reliability [like the error-spread score (Christensen et al. 2015)] depend on the third or higher
 421 order moments of the forecast probability distribution. Since the StocSIPS forecast is Gaussian by definition, the ESS used
 422 here (Eq. (12)) gives enough information assuming the near Gaussianity of the observational probability distribution.

423 In Fig. 8 we show maps of the ESS of StocSIPS for horizon from 1 to 4 months. Notice that, from Eq. (A13), $\sigma_{\text{ensemble}} =$
 424 $\text{RMSE}_{\text{nat}}^{\text{theory}}$ is a function of the forecast horizon and the location – following the spatial distribution of the model parameters
 425 σ_T and H –, but for all pixels the ESSs are very close to 1, except for the tropical ocean where it tends to be “overdispersive”
 426 ($\text{ESS} > 1$). The average values for the globe with one standard deviation are shown in brackets in each map label. They increase
 427 monotonically from 0.96 ± 0.05 for $k = 1$ month, 0.98 ± 0.03 for $k = 2$ months, 0.98 ± 0.04 for $k = 3$ months, $1.00 \pm$
 428 0.06 for $k = 4$ months, up to 1.09 ± 0.21 for $k = 12$ months (only the first four values are included in the maps). From Eq.
 429 (11), it can be shown that for a system with perfect reliability where $\text{ESS} = 1$, the CRPS takes its minimum value $\text{CRPS}_{\text{min}} =$
 430 $\text{RMSE}/\sqrt{\pi}$. For any other case when we have an “overconfident” ($\text{ESS} < 1$) or an “overdispersive” ($\text{ESS} > 1$) system, $\text{CRPS} >$
 431 $\text{RMSE}/\sqrt{\pi}$. In conclusion, StocSIPS is a nearly perfectly reliable system (except for the tropical ocean) without need of
 432 recalibration of the forecast probability distribution.

433 3.2 Hindcast verification

434 3.2.1 Monthly and 3-month average predictions

435 The results presented in Sect. 3.1 confirm the validity of the stochastic model on forecasting the natural temperature variability.
 436 In this section, we show the verification scores for the forecast of the raw (undetrended) anomalies including the forecast of
 437 the CO_2eq deterministic trend. All the scores were computed following the definitions shown in Appendix B.

438 Given the smooth variation of the CO_2eq concentration at monthly scales, we can use simple extrapolation in Eq. (2) to obtain
 439 the predictor $\hat{T}_{\text{anth}}(t + k)$ from the knowledge of the CO_2eq concentration path up to time t . As the function $T_{\text{anth}}(t)$ is almost
 440 linear in a k -vicinity of any t , the error of projecting the anthropogenic component is negligible compared to the error of the

441 natural variability. In fact, as we assume the same global CO₂eq forcing affecting all locations, the error of predicting the
 442 anthropogenic trend for a given k , is proportional to the sensitivity map shown in Fig. 4c. It was found that this error is lower
 443 than 2% of the RMSE of the natural variability for all locations and for all horizons. In any case, the projection of the trend
 444 was still included in the following verification results.

445 **Normalized Root Mean Square Error (NRMSE)**

446 Figure 9 shows maps of the NRMSE for horizons $k = 1, 2$ and 3 months (panels (a), (b) and (c), respectively) and for the
 447 seasonal forecast (including all seasons, average for $k = 1 - 3$ months) in panel (d). The values in brackets in the figure labels
 448 are the NRMSE globally area averaged over the grid points (see Eq. (B9)). In general, the skill of the forecasts is larger over
 449 ocean than over land, with the lower values of NRMSE attained over the tropical ocean. This corresponds to the distribution of
 450 H shown in Figs. 2 and 4a.

451 Since small NRMSE implies large skill, according to the global-averaged NRMSE, the seasonal skill is larger than that of any
 452 of the first three individual monthly forecasts. This is possible because although the horizon is further in the future, the seasonal

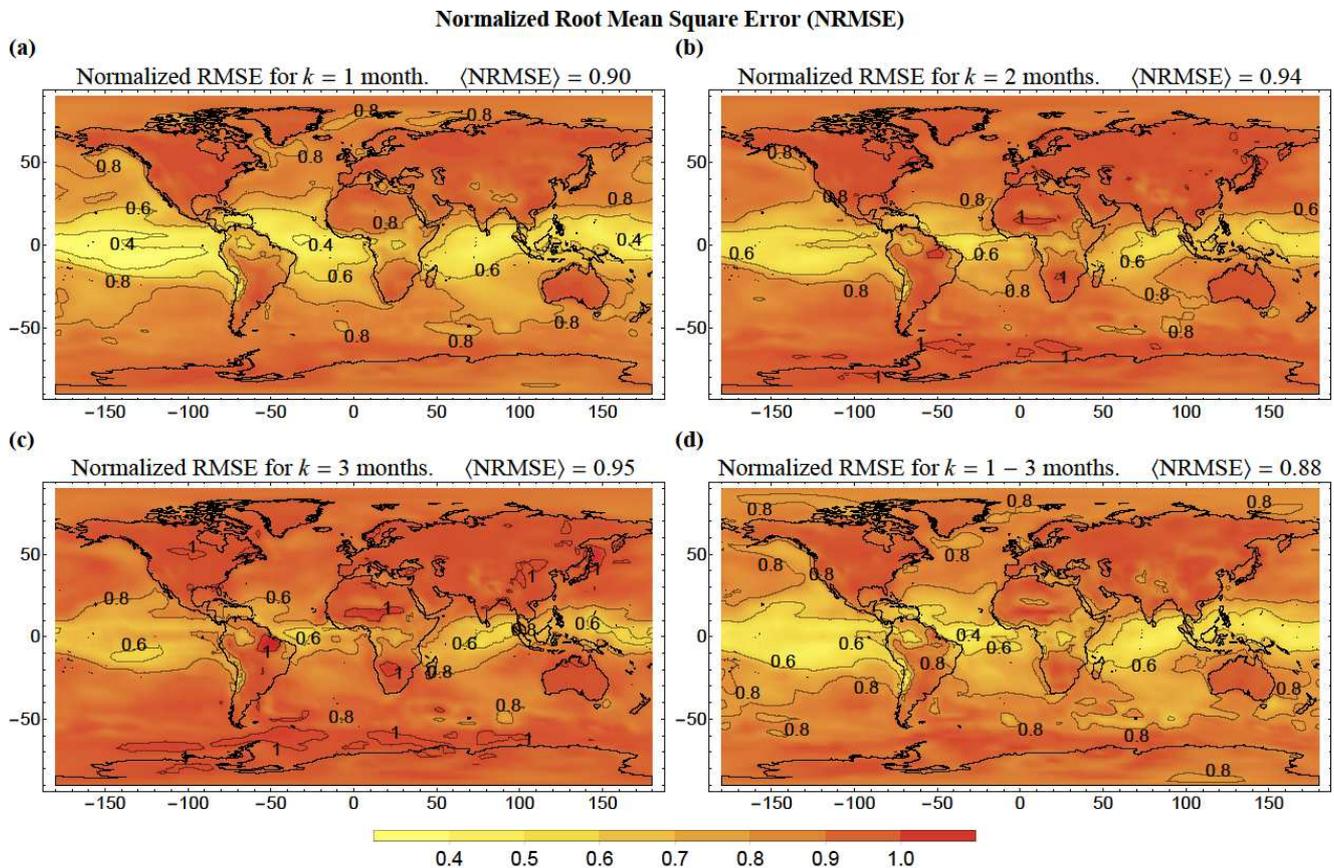


Fig. 9 Normalized root mean square error NRMSE for: (a) $k = 1$ month, (b) $k = 2$ months, (c) $k = 3$ months and (d) for the all-seasons mean (average for $k = 1 - 3$ months). The values in brackets in the figure labels represent the areal mean of global NRMSE.

453 forecast is for a longer (3 month) average. For scaling processes, the two effects exactly compensate. For the prediction of the
 454 natural variability component using fGn, the skill on predicting the next month using monthly averaged data is the same as the
 455 skill on predicting the next season using 3-month averaged data. This is reflected in Eq. (A13), where k is in units of τ , which
 456 is the resolution (smallest sampling temporal scale) of the data. The similarity between the average values in the captions of
 457 panels (a) and (d) of Figs. 9-11 confirms this consequence of the scaling.

458 The values in Fig. 9a for the forecast of the raw anomalies are lower than those shown in Fig. 5b for the natural variability
 459 because, while the RMSE of both are almost the same (we can neglect the error on projecting the anthropogenic trend), the
 460 normalization factor (standard deviation of the respective anomalies) is larger for the undetrended anomalies.

461 **Mean Square Skill Score (MSSS)**

462 To compute the MSSS for the raw anomalies, the MSE of the reference climatology forecasts (forecast produced using only the
 463 annual cycle signal without removing the anthropogenic variation) is in this case:

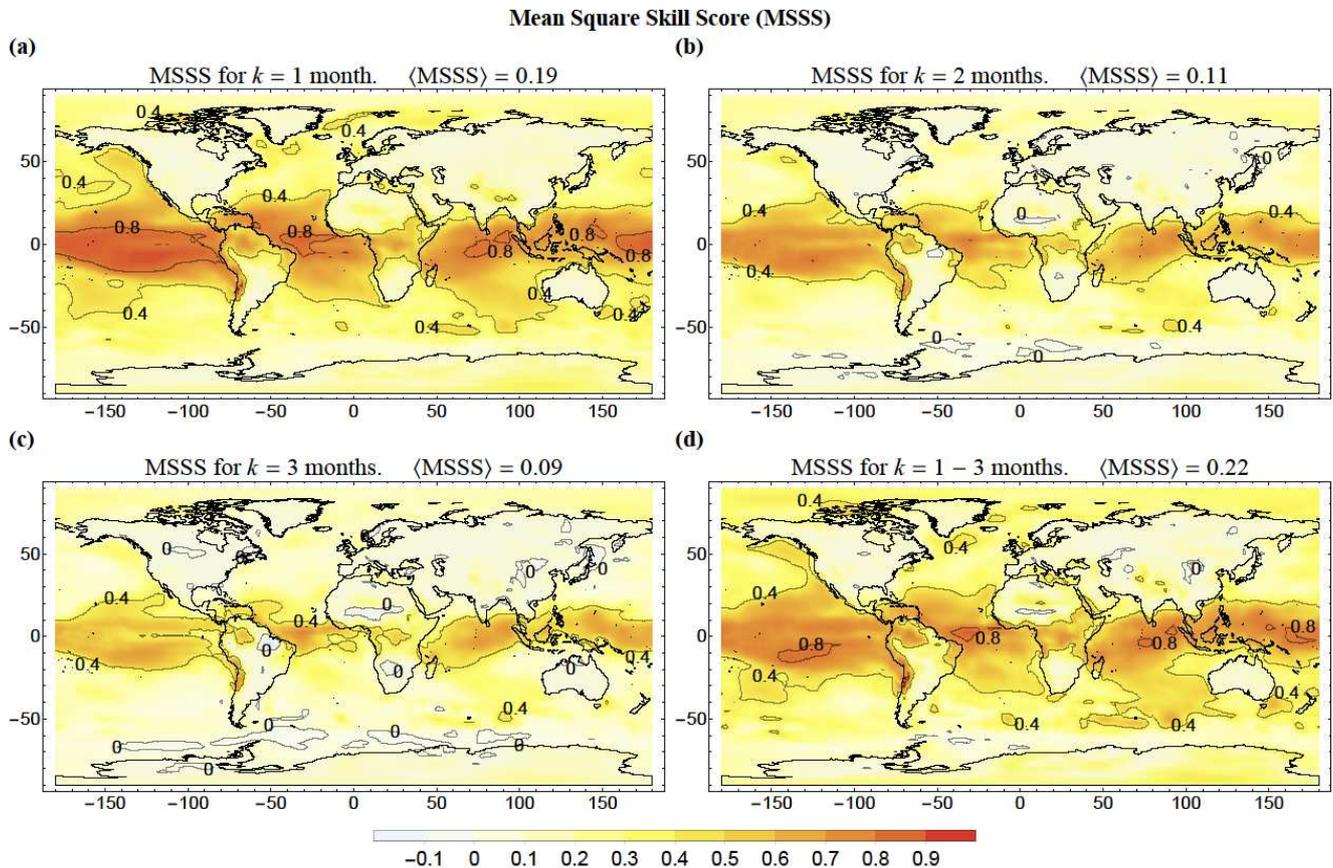


Fig. 10 Mean square skill score (MSSS) for: (a) $k = 1$ month, (b) $k = 2$ months, (c) $k = 3$ months and (d) for the all-seasons mean (average for $k = 1 - 3$ months). The values in brackets in the figure labels represent the areal mean of global MSSS.

464

$$\text{MSE}_c = \left(\frac{N}{N-1} \right)^2 SD_{\text{anom}}^2, \quad (13)$$

465 where SD_{anom}^2 is the variance of the anomalies series without removing the anthropogenic component:

$$SD_{\text{anom}}^2 = \overline{T_{\text{anom}}^2} = \overline{(T_{\text{anth}} + T_{\text{nat}})^2} = \overline{T_{\text{anth}}^2} + SD_T^2 \quad (14)$$

467 (assuming that the natural and anthropogenic variabilities are independent).

468 Because $SD_{\text{anom}}^2 > SD_T^2$ and the MSE of the forecast of the raw and the detrended anomalies are almost equal, then from Eq.
469 (8) we obtain that the MSSS for the undetrended series forecast is larger than for the natural variability.

470 Maps of MSSS, corresponding to those shown in Fig. 9, are shown in Fig. 10. The difference in skill between ocean and land
471 is more evident in these maps. In many places over land, the MSSS is close to zero, meaning that most of the skill comes from
472 the projection of the anthropogenic trend. The global averages shown in brackets in the map labels are computed following
473 the guidelines of the WMO (WMO 2010a). Note that the maps and the average values shown in Figs. 9 and 10 are related as
474 $\text{MSSS} \approx 1 - \text{NRMSE}^2$ if the reference forecast for the MSSS is the climatological annual cycle.

Temporal Correlation Coefficient (TCC)

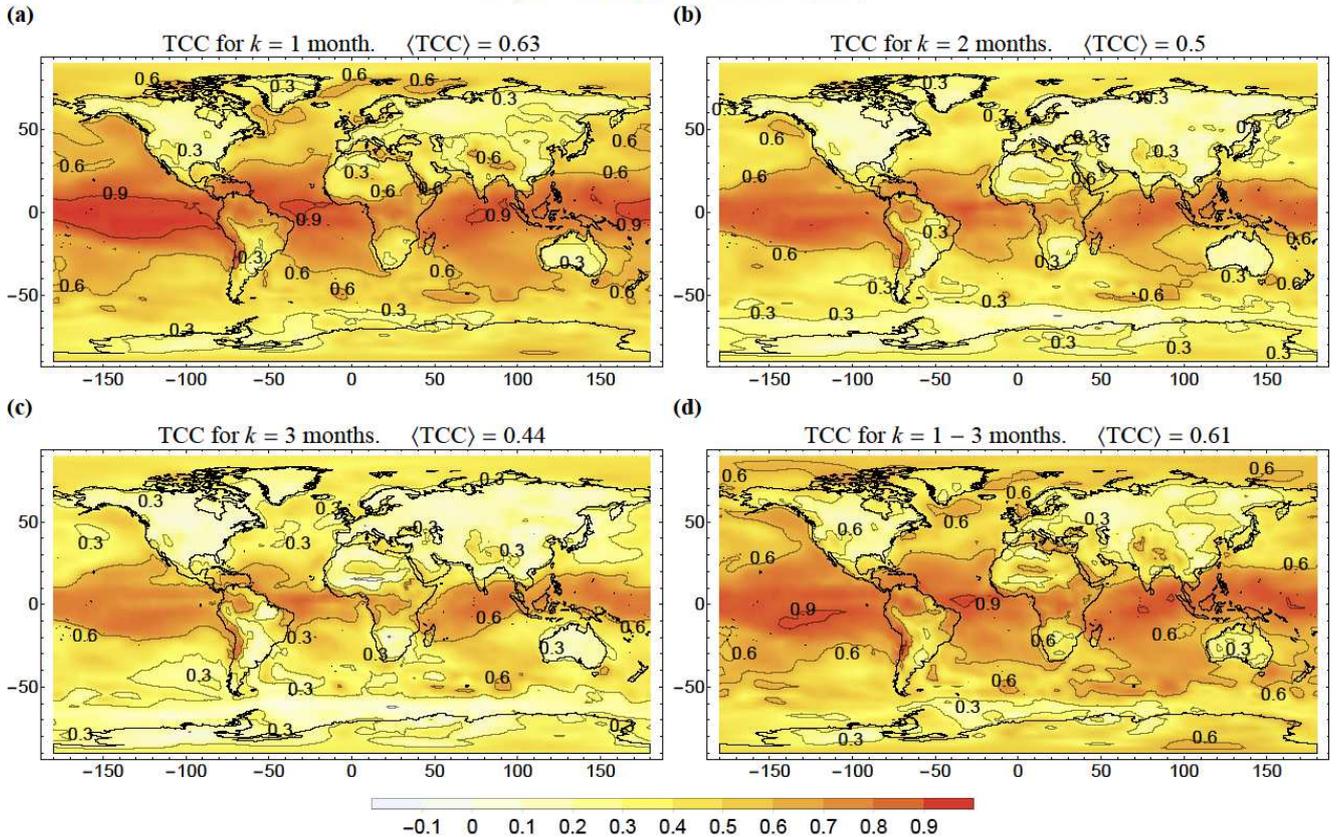


Fig. 11 Anomaly correlation coefficient (TCC) for: (a) $k = 1$ month, (b) $k = 2$ months, (c) $k = 3$ months and (d) for the all-seasons mean (average for $k = 1 - 3$ months). The values in brackets in the figure labels represent the areal mean of global TCC.

475 **Temporal Correlation Coefficient (TCC)**

476 Similarly to the MSSS, if the TCC is obtained for the undetrended anomalies (with only the annual cycle, but not the
477 anthropogenic trend removed), then higher values are often obtained compared to the TCC for the natural variability because
478 of the extra correlation associated to the trend. For most of the long-term forecasts reported in the literature, the increasing
479 trend related to anthropogenic warming is not removed and only the annual cycle is considered to obtain the anomalies used
480 for verification.

481 In Fig. 11, we show maps of the TCC for the prediction of the raw anomalies. The number in brackets in the caption of each
482 plot indicates the area-averaged over the globe of the grid-point correlation coefficients. The area average was computed taking
483 the Fisher Z-transform of the correlations following Eq. (B12) (Fisher 1915; WMO 2010b). The StocSIPS predictions over the
484 ocean are highly correlated with the observations and the highest correlations are in the tropical regions. Over land, although
485 the skill is poorer (using the correlation coefficient), it is still significantly high for the forecast of the first three months. The
486 TCC of the prediction is positive almost everywhere and, compared to the NRMSE or the MSSS, it shows significantly larger
487 skill. This “extra” skill shown in the correlations for the raw anomalies comes from the presence of the anthropogenic signal.

488 **3.2.2 Global Averages**

489 To summarize, in Fig 12 we show graphs of the area-averaged NRMSE, MSSS and TCC for the monthly and the 3-month
490 average forecasts as a function of the forecast horizon. In all the graphs, the red lines with circles correspond to the average
491 considering the grid points for the whole planet, the blue lines with open squares are for places over the ocean and the green
492 lines with triangles are for grid points over land. The corresponding dashed lines of the same colours represent the respective
493 scores obtained if only the anthropogenic trend is forecast. In all cases, the reference forecast is the climatological annual
494 cycle. Attending to the average values, we can conclude that the skill over ocean is always greater than over land, with the
495 global skill in between the two.

496 As we mentioned previously, for a perfectly scaling process, the 3-month average forecasts for $k = 1 - 3$ months would have
497 the same skill as the monthly forecast for $k = 1$ month. In the same way, the seasonal for $k = 4 - 6$ months would correspond
498 to the monthly for $k = 2$ months, for $k = 7 - 9$ months to $k = 3$ months and for $k = 10 - 12$ months to $k = 4$ months. A
499 comparison between panels (a) and (d) and (b) and (e) in Fig. 12, show that this is reasonably well satisfied for the NRMSE
500 and MSSS, respectively. Of course, the actual comparison should be made for the forecast of the natural temperature variability,
501 which is the true scaling process. The results shown in Fig. 12 are for the raw anomalies which include the anthropogenic
502 trend, that breaks the scale invariance of the fluctuations.

503 The difference between the curves and the dashed horizontal lines (showing the skill if only the anthropogenic trend is forecast)
504 corresponds to the skill on forecasting the natural variability using the fGn model. While the forecast of the natural variability
505 is reasonably skillful for $k \leq 6$ months, for horizons larger than 6 months, most of the overall skill comes from projecting the
506 anthropogenic trend. Finally, note that the global and ocean averages vary monotonically with k , but the land averages show
507 some oscillation that indicates a seasonality effect in the forecasts. This seasonality is analyzed in the next section.

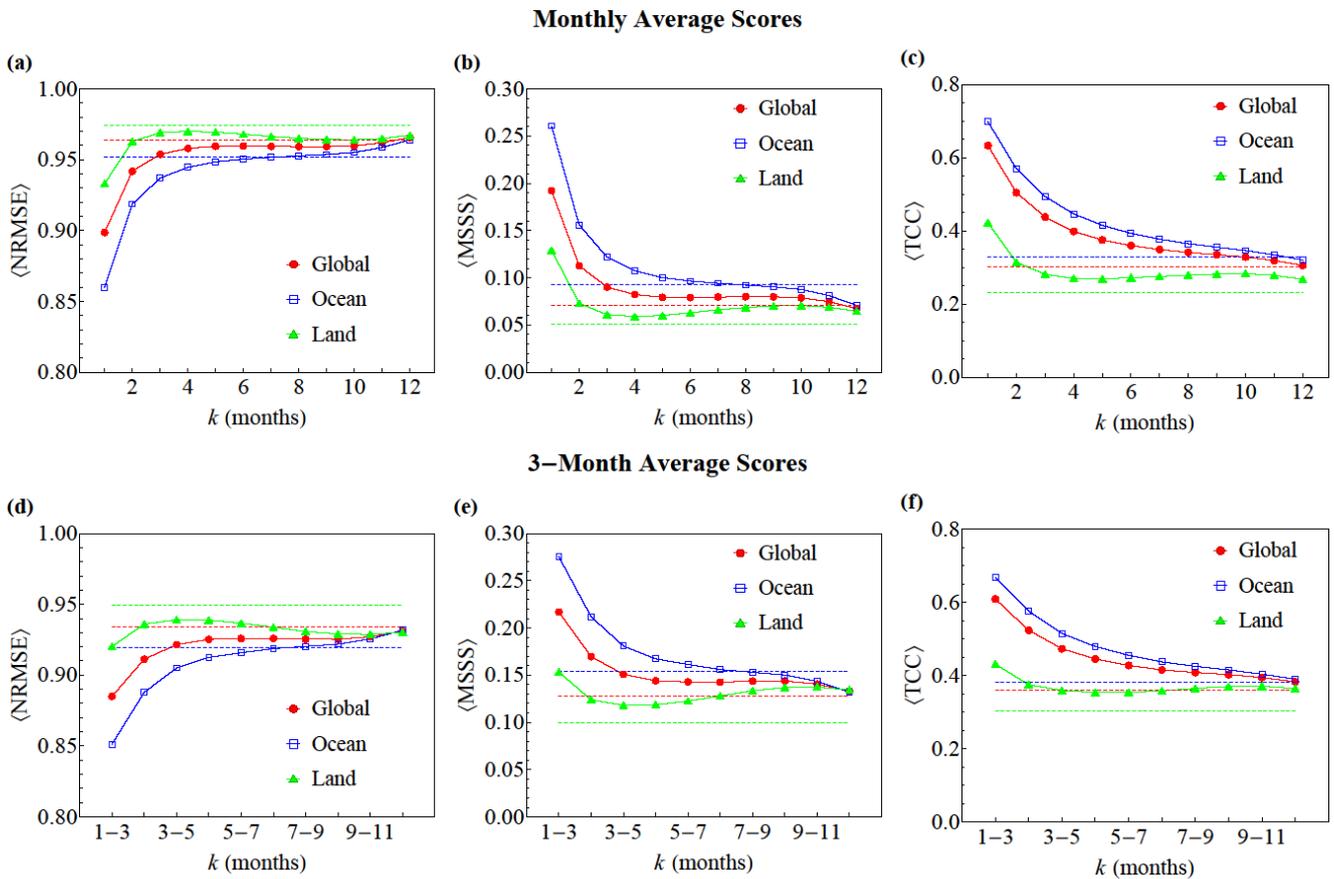


Fig. 12 Graphs of the area-averaged NRMSE, MSSS and ACC for the monthly (panels (a), (b) and (c)) and the 3-month average (panels (d), (e) and (f)) forecasts as a function of the forecast horizon. In all the graphs, the red lines with circles correspond to the average considering the grid points for the whole planet, the blue lines with open squares are for places over the ocean and the green lines with triangles are for grid points over land. The corresponding dashed lines of the same colours represent the respective scores obtained if only the anthropogenic trend is forecast.

508 3.2.3 Multiplicative seasonality

509 The results shown in panel (d) of Figs. 9-11, were obtained for the seasonal forecast without distinguishing specific seasons.
 510 In fact, StocSIPS assumes that each month has the same anomaly statistics. It is actually this month-to-month correlation that
 511 is exploited as a source of predictability in the stochastic model. Nevertheless, there is always an intrinsic multiplicative
 512 seasonality in the data that is impossible to completely remove without affecting the scaling behaviour. This seasonal
 513 interannual variability is shown in Fig. 13, where the standard deviation of the 3-month averaged anomalies is shown for each
 514 conventional season: (a) December to February (DJF), (b) March to May (MAM), (c) June to August (JJA) and (d) September
 515 to November (SON). The difference in the variability between the spring and the fall seasons (panels (b) and (d)) is low. In
 516 comparison, the interannual variability over the land area in the northern hemisphere is larger during the boreal winter (DJF)
 517 and lower during the summer (JJA).

Seasonal Interannual Variability

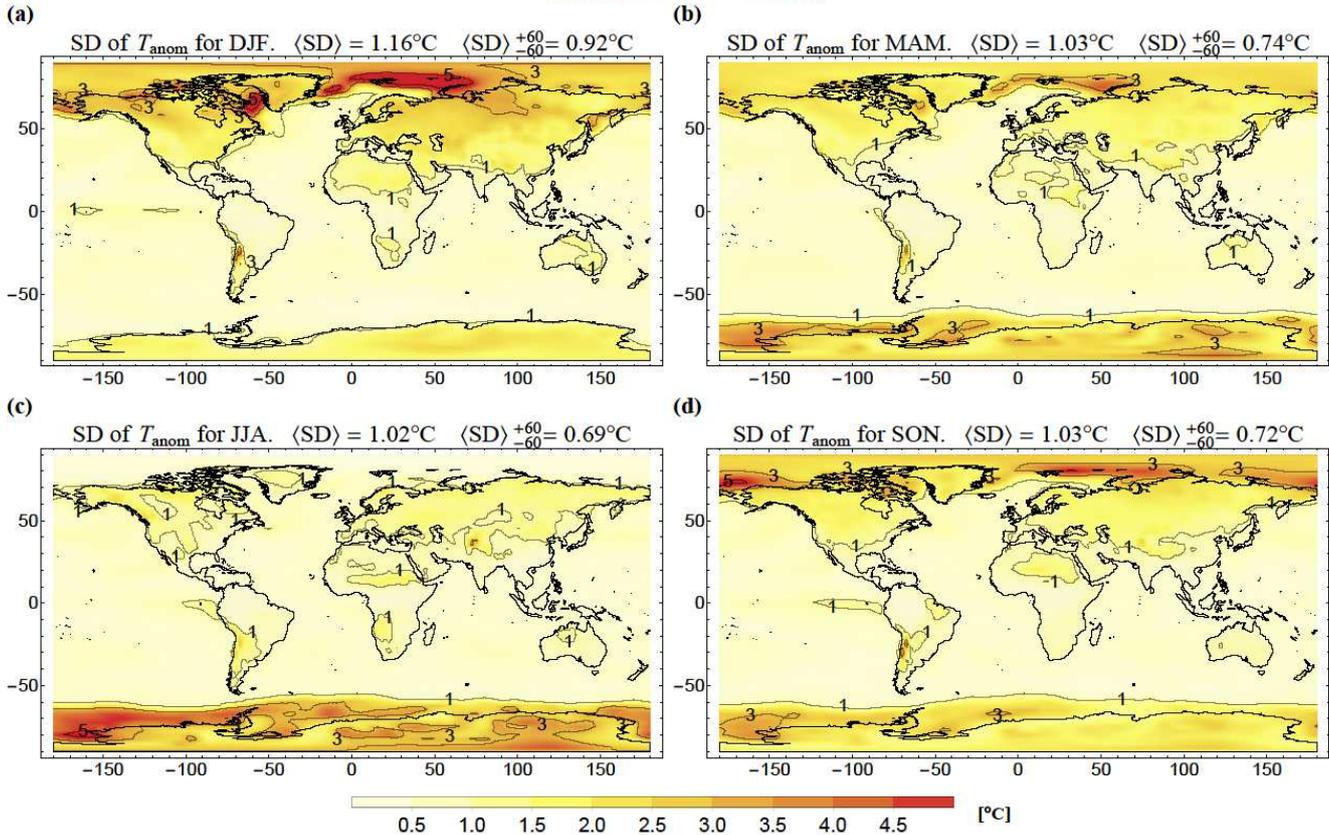


Fig. 13 Interannual standard deviation (SD) of the temperature anomalies for the conventional seasons: (a) DJF, (b) MAM, (c) JJA and (d) SON. The values in brackets in the figure labels represent the areal mean of global standard deviation and the areal mean excluding the poles (between 60°S and 60°N).

518 The largest seasonality is observed in the polar regions, where the winter temperatures are much more intermittent compared
 519 to the summer values. Conversely, during the summer the standard deviation of the anomalies in the poles is much lower
 520 compared to the other seasons. The values in brackets in the figure labels represent the areal mean of global standard deviation,
 521 $\langle SD \rangle$, and the areal mean excluding the poles (between 60°S and 60°N), $\langle SD \rangle_{-60}^{+60}$. Notice that the polar regions contribute
 522 substantially to the interannual variability and also, that the boreal winter season is in general significantly more variable than
 523 the other seasons (the $\langle SD \rangle$ goes from 1.16°C for DJF to roughly 1.03°C for the others). A possible explanation for this
 524 seasonality is that, when removing the annual cycle and the trend associated with the anthropogenic warming, we assumed
 525 that both were statistically independent. This is not so true for the polar region. While for the rest of the planet the anthropogenic
 526 temperature response increases uniformly for every month following the increasing CO₂ concentrations, in the poles during
 527 the summer, the temperature is tied to the freezing point of water. This is a shortcoming of the model that could be considered
 528 to improve future versions of StocSIPS.

529 3.2.4 Preliminary comparison with GCMs' seasonal predictions

530 In the previous sections, we validated StocSIPS as a good model for describing the monthly surface temperature field and we
531 assessed its skill by computing monthly and 3-month average scores, without distinguishing specific seasons. To account for
532 the effects of the multiplicative seasonality on the predictions, we can stratify the observations and the forecasts series to show
533 dependencies with the targeted season and the forecast horizon. Usually, this is the kind of forecast published by several major
534 operational centers for seasonal prediction. In this section we show the skill scores obtained for stratified data and we make a
535 preliminary comparison with other models' skill to assess the relative advantages and shortcomings of StocSIPS.

536 Our purpose in this paper is not to make an exhaustive and detailed comparison with other long-term prediction models' results.
537 This detailed comparison and also the combination of StocSIPS with conventional numerical models to produce merged
538 forecasts is the subject of a future paper currently in preparation. Those results are too extensive to include them in the present
539 paper, so we limited ourselves to compare with already published skill scores from other models. For this purpose, we selected
540 the multi-model ensemble (MME) predictions recently published by (Kim et al. 2020). An important aspect is that Kim et al.
541 offer a detailed description of the scores and the methods used, which we try to closely follow here to guarantee reproducibility.
542 The definition of these metric are given in Appendix B, following the guidelines of the WMO Standardized verification system
543 for long-range forecasts (SVS-LRF) (WMO 2010b, a).

544 In (Kim et al. 2020), the authors asses different MME combination methods for seasonal prediction using hindcast datasets of
545 six models from five Global Producing Centers (GPCs) for long-range forecasts (LRFs) designated by the WMO (Graham et
546 al. 2011). The six models included in their analysis cover 27 years of common hindcast period from 1983 to 2009. The selected
547 GPCs were: Melbourne, Montreal (two models), Moscow, Seoul and Tokyo. References and details of the individual models
548 can be found in (Kim et al. 2020). The authors study seven experimental deterministic MME methods to merge the six seasonal
549 forecast systems: simple composite method (SCM), simple linear regression (SLR), multiple linear regression (MLR), best
550 selection anomaly (BSA), multilayer perceptron (MLP), radial basis function (RBF) and genetic algorithm (GA). Their
551 reported scores for 2-m temperature were obtained for 1-month lead retrospective forecasts in a grid with a resolution of 2.5°
552 in both longitude and latitude. To produce the figures in this section, we used and adapted some of the figures from (Kim et
553 al. 2020) (including supporting information).

554 **Mean Square Skill Score (MSSS)**

555 For a better comparison with Kim et al. results, all the seasonal scores for StocSIPS were obtained for observational and
556 forecast seasonal anomalies calculated as departures from the climatology in the leave-one-out cross-validation scheme for the
557 period 1983-2009. In Fig. 14, we show maps of the MSSS of StocSIPS for: (a) DJF, (b) MAM, (c) JJA and (d) SON. In all
558 cases, the forecasts used data up to the beginning of each respective season (average for $k = 1 - 3$ months), i.e. including
559 November for DJF, up to February for MAM and so on. The values in brackets in the figure labels represent the globally

560 averaged score, $\langle \text{MSSS} \rangle$, (see Eq. (B10)). In panels (e) and (f) we reproduce maps of MSSS for DJF and JJA, respectively, from
 561 Figs. S1 and S2 of (Kim et al. 2020) (supporting information) for their best MME combination method (GA).

Seasonal Mean Square Skill Score (MSSS)

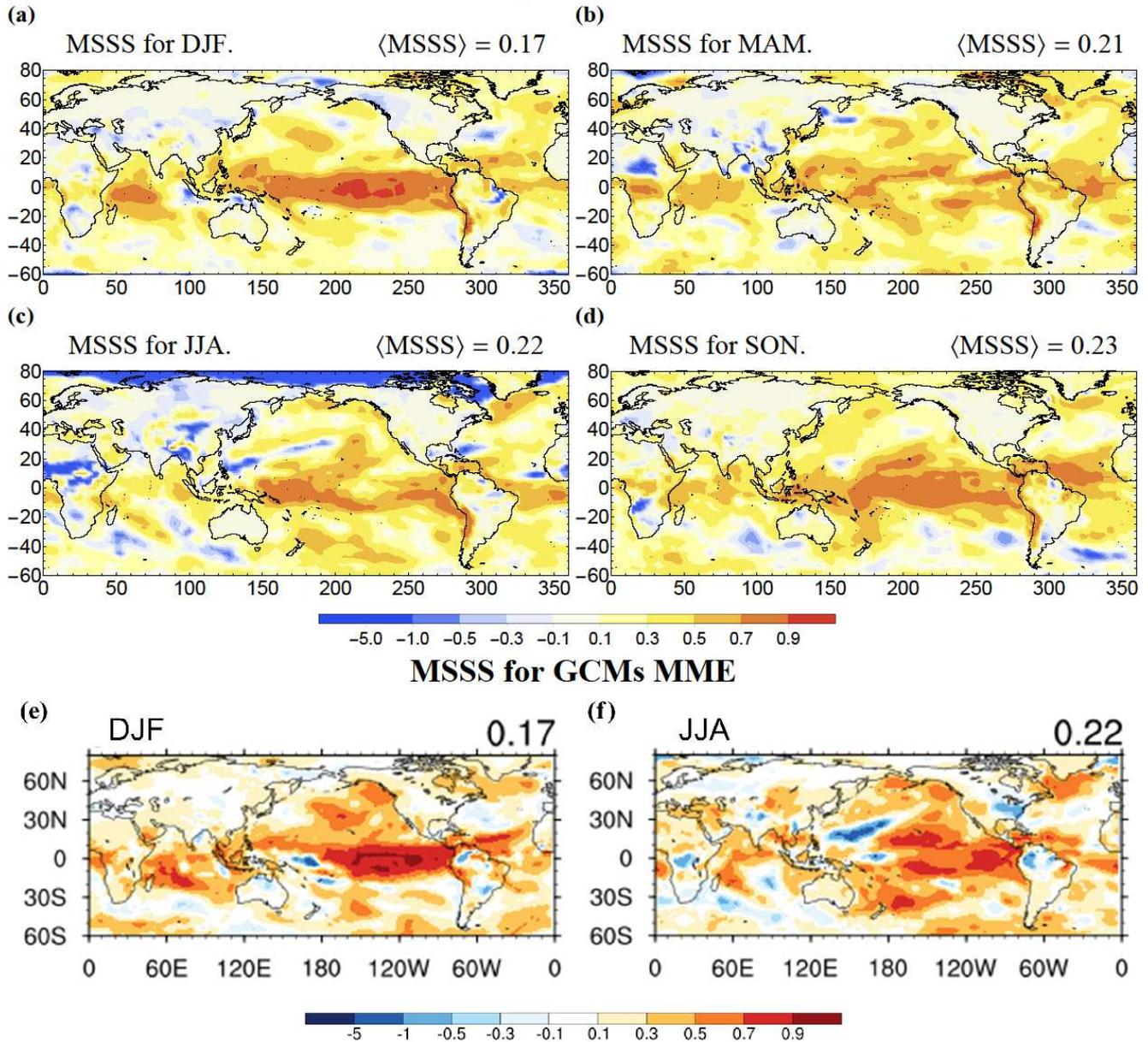
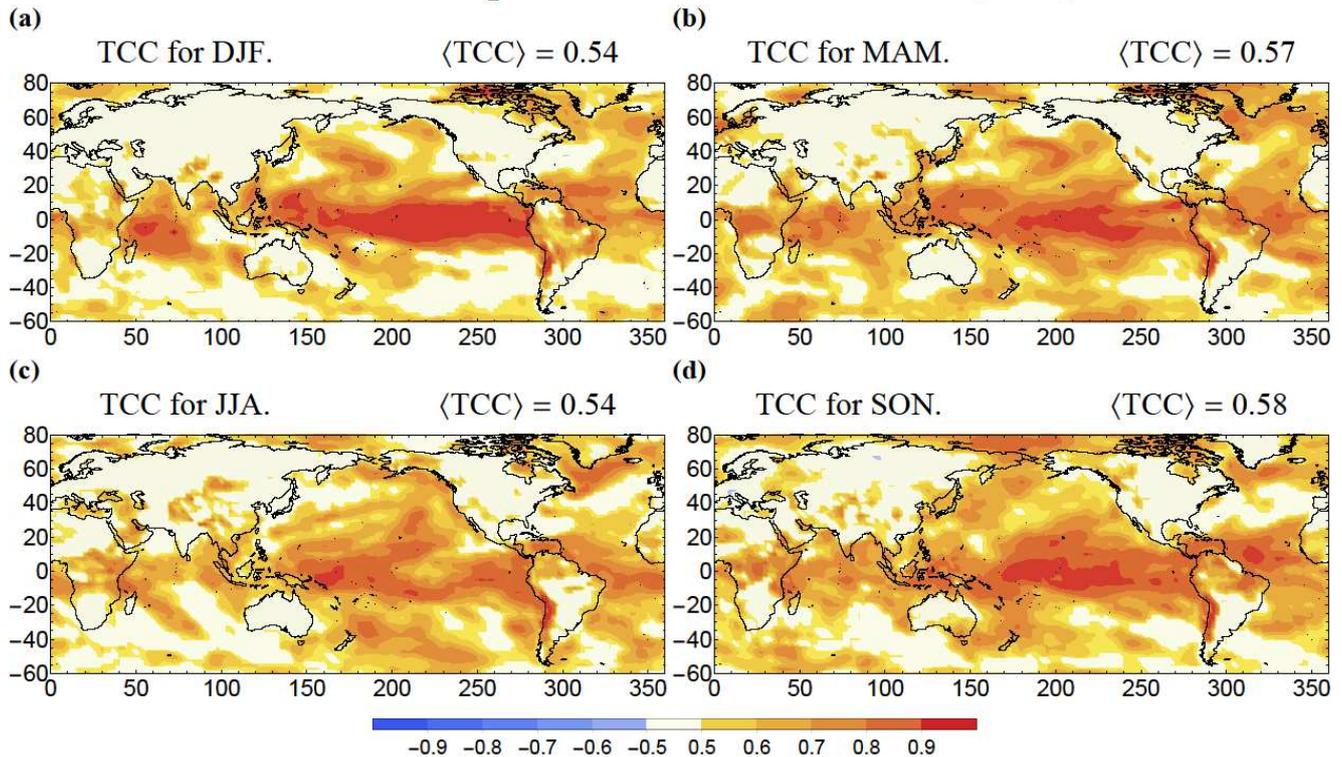


Fig. 14 MSSS for: (a) DJF, (b) MAM, (c) JJA and (d) SON. In all cases, the forecasts used data up to the beginning of each respective season (average for $k = 1 - 3$ months). The values in brackets in the figure labels represents the globally averaged MSSS (see Eq. (B10)). The maps shown in panels (e) and (f) for the GCMs MME prediction of DJF and JJA, respectively, were reproduced from Figs. S1 and S2 of (Kim et al. 2020) (supporting information) for their best MME combination method (GA). This is an open access article distributed under the terms of the Creative Commons CC BY license, which permits unrestricted use, distribution, and reproduction in any medium.

Seasonal Temporal Correlation Coefficient (TCC)



TCC for GCMs MME

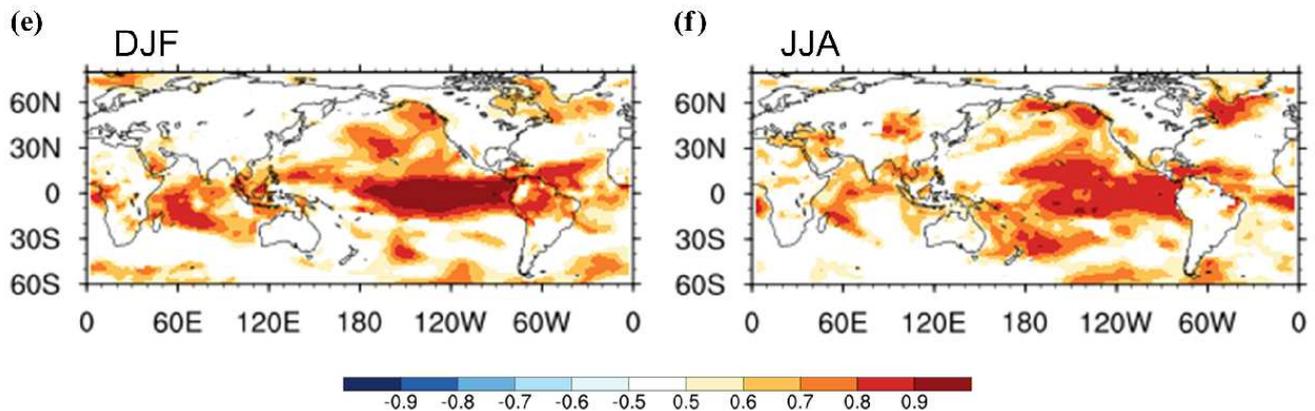


Fig. 15 TCC for: (a) DJF, (b) MAM, (c) JJA and (d) SON. In all cases, the forecasts used data up to the beginning of each respective season (average for $k = 1 - 3$ months). The shaded areas indicate the regions over the 5% significance level using two-tailed student's t-test. The values in brackets in the figure labels represent the globally averaged score, $\langle \text{TCC} \rangle$, computed using Eq. (B12). The maps shown in panels (e) and (f) for the GCMs MME prediction of DJF and JJA, respectively, were reproduced from Figs. S5 and S6 of (Kim et al. 2020) (supporting information) for their best MME combination method (GA). This is an open access article distributed under the terms of the Creative Commons CC BY license, which permits unrestricted use, distribution, and reproduction in any medium.

563 for the individual seasons in general show better skill over the ocean than over land. The MSSS values are particularly high in
564 the tropical region with the highest values obtained in the equatorial Pacific for DJF. Similar results were obtained for the
565 GCM forecasts shown in Fig. 14 (e) and (f). The globally averaged scores (shown in the top right corner of each plot), are
566 identical for StocSIPS and the MME results: 0.17 and 0.22 for DJF and JJA, respectively. The negative values of MSSS for
567 StocSIPS near the north pole for JJA are associated to the multiplicative seasonality effect described in Sect. 3.2.3.

568 **Temporal Correlation Coefficient (TCC)**

569 Similar to Fig. 14 for the MSSS, in Fig. 15 we show maps of the TCC of StocSIPS for: (a) DJF, (b) MAM, (c) JJA (d) SON and
570 the best GCMs MME combination (GA) from (Kim et al. 2020) in panels (e) and (f) for DJF and JJA, respectively. The shaded
571 areas indicate the regions over the 5% significance level using two-tailed student's t-test. The values in brackets in the figure
572 labels represent the globally averaged score, $\langle TCC \rangle$, computed using Eq. (B12). As before, the highest correlation values are
573 achieved in tropical regions. Considering the average scores, there is no significant reduction in the TCC for DJF compared to
574 JJA. There are also no considerably low values near the north pole for JJA. Compared to the MSSS, the multiplicative
575 seasonality effects are less reflected in the TCC, since the latter is a measure of the skill in predicting the phase (sign), so less
576 dependent on the variability of the anomaly magnitudes.

577 **Anomaly Pattern Correlation Coefficient (TCC)**

578 The temporal evolution of the forecast skill can be assessed using the anomaly pattern correlation coefficient (ACC), which is
579 the spatial correlation for any given date between the observational and forecast anomalies (see Eq. (B7)). This shows how
580 well the model reproduces the temperature anomaly distribution around the globe for any given season. Figure 16 shows the
581 evolution of the ACC for StocSIPS (black line with solid circles) and for each of the seven MME combination methods studied
582 by Kim et al. (colored lines with markers) in the 27-year verification period 1983-2009. Graphs for each season are shown
583 independently: (a) MAM, (b) JJA, (c) SON and (d) DJF. This figure was adapted from Fig. 3 in (Kim et al. 2020) to include
584 the StocSIPS scores. In the original figure, the authors also show the absolute value of the El Niño 3.4 index (black line without
585 markers) to study the dependence of the ensemble predictions with the El Niño phase. The main conclusion is that the GCM
586 MMEs perform better during ENSO events than during non-ENSO events for all seasons. A similar behaviour was not found
587 for the case of StocSIPS, where the performance based on the ACC varies independently of the ENSO phase. The average
588 scores for the POV (see Eq. (B13)) are shown in the right panels for each of the respective seasons. Comparing these values,
589 we can see that StocSIPS has better overall skill than most of the GCM MME combinations for all seasons. Only for JJA, the
590 StocSIPS score is lower than the best three MME (SCM, SLR and GA). For the rest of the seasons, its average score is almost
591 equal (or slightly larger) than the best MME (using GA or SCM).

592 **Globally averaged TCC and RMSE**

593 Comparisons for globally averaged TCC and RMSE (see Eqs. (B9) and (B12)) are shown in Fig. 17 (a) and (b), respectively,
594 for each season. The bars are for the MME combination methods in (Kim et al. 2020), together with the mean of single model
595 skills (MSMS). The scores for StocSIPS were included as horizontal lines with the same color code as the bars for each
596 respective season. The dashed black line indicates that the estimated TCC is statistically significant at the 5% level using the

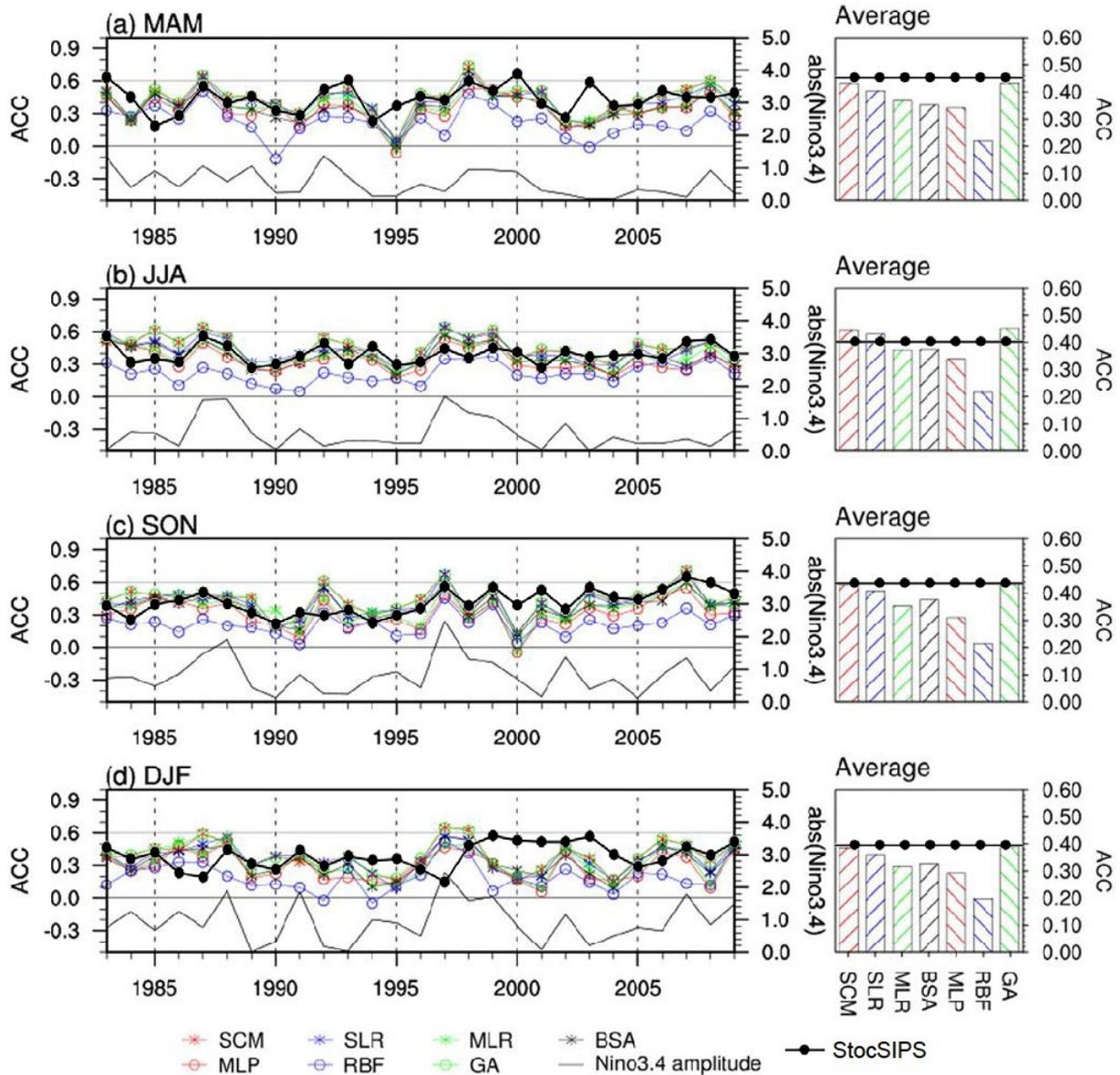


Fig. 16 ACC for StocSIPS (black line with solid circles) and for each of the seven MME combination methods studied by Kim et al. (colored lines with markers) in the 27-year verification period 1983-2009 for: (a) MAM, (b) JJA, (c) SON and (d) DJF. The average scores for the POV (see Eq. (B13)) are shown in the right panels for each of the respective seasons. The absolute value of the El Niño 3.4 index (black line without markers) is also shown. This figure was adapted from Fig. 3 in (Kim et al. 2020) to include the StocSIPS scores. This is an open access article distributed under the terms of the Creative Commons CC BY license, which permits unrestricted use, distribution, and reproduction in any medium.

597 one-tailed Student's t test. The GA methods shows the best performance, although it is very close to the SCM with equal
 598 weights for each model. Most MME predictions show higher skill than the corresponding MSMS for all four seasons, although
 599 sometimes (like the TCC for MLP), the MME combination does not improve over the single model predictions. In all cases,

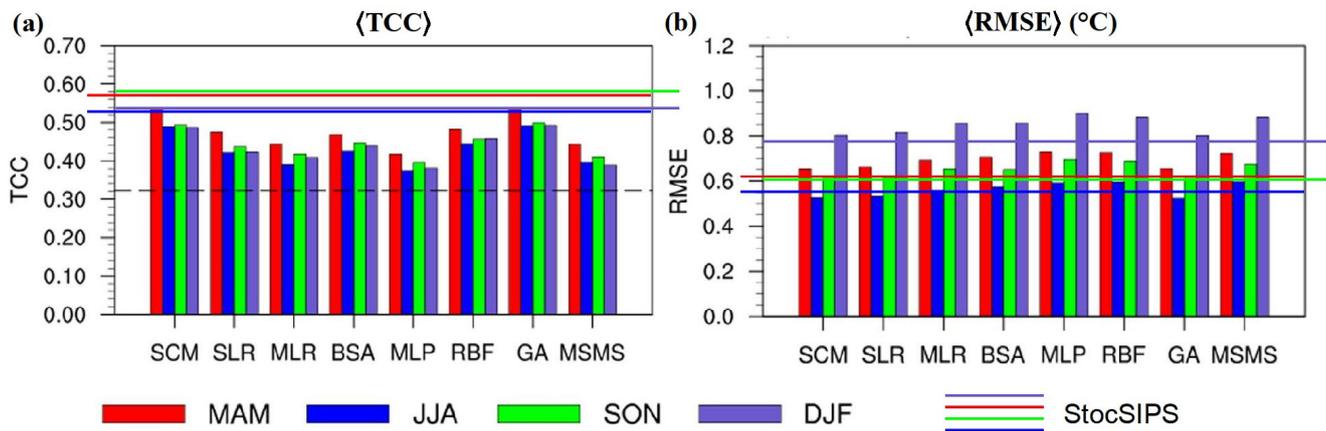


Fig. 17 Globally averaged TCC (a) and RMSE (b) (Eqs. (B12) and (B9), respectively) for MAM (red), JJA (blue), SON (green) and DJF (purple) for the period 1983-2009. The bars are for the MME combination methods in (Kim et al. 2020), together with the mean of single model skills (MSMS). The scores for StocSIPS were included as horizontal lines with the same color code for each respective season. The dashed black line indicates that the estimated TCC is statistically significant at the 5% level using the one-tailed Student's t test. This figure was adapted from Figs. 5 and 6 in (Kim et al. 2020) to include the StocSIPS scores. This is an open access article distributed under the terms of the Creative Commons CC BY license, which permits unrestricted use, distribution, and reproduction in any medium.

600 the TCC of the StocSIPS forecasts is larger than the best GCM MME. Similarly, the StocSIPS RMSE is lower than most of the
 601 MME combinations for all seasons. Only the SCM, GA and SLR show lower errors than StocSIPS for JJA predictions. The
 602 globally averaged TCC does not show a large seasonal variation, but there is still a reduction in skill for JJA and DJF associated
 603 to the high variability in the poles discussed in Sect. 3.2.3. This multiplicative seasonality effect is clear in the average RMSE,
 604 which follows the average SD values in the caption of Fig. 13.

605 For an overall comparison, in Fig. 18 we show a plot of the 4-season-averaged RMSE vs. TCC for the six individual models
 606 used in (Kim et al. 2020) (red crosses) and the six MME combinations (letters). The StocSIPS scores were included as a blue
 607 asterisk. For the GCMs, the GA method has the best performance – very close to the SCM – with the highest TCC (0.51) and
 608 the lowest RMSE (0.64). The StocSIPS forecasts have similar RMSE (0.64), but better average TCC (0.55).

609 4 Summary and discussion

610 In this paper we applied the Stochastic Seasonal to Interannual Prediction System (StocSIPS) to the monthly and seasonal
 611 prediction of the surface temperature with a $2.5^\circ \times 2.5^\circ$ spatial resolution. The theory and the basis of the numerical methods
 612 used in StocSIPS were previously presented and applied to the forecast of globally averaged temperature in (Del Rio Amador
 613 and Lovejoy 2019). StocSIPS is based on two statistical properties of the macroweather regime: the near Gaussianity of
 614 temperature fluctuations and the temporal scaling symmetry of the natural variability. The model is a high-frequency
 615 approximation to the Fractional Energy Balance Equation (FEBE), a fractional generalization of the usual EBE.

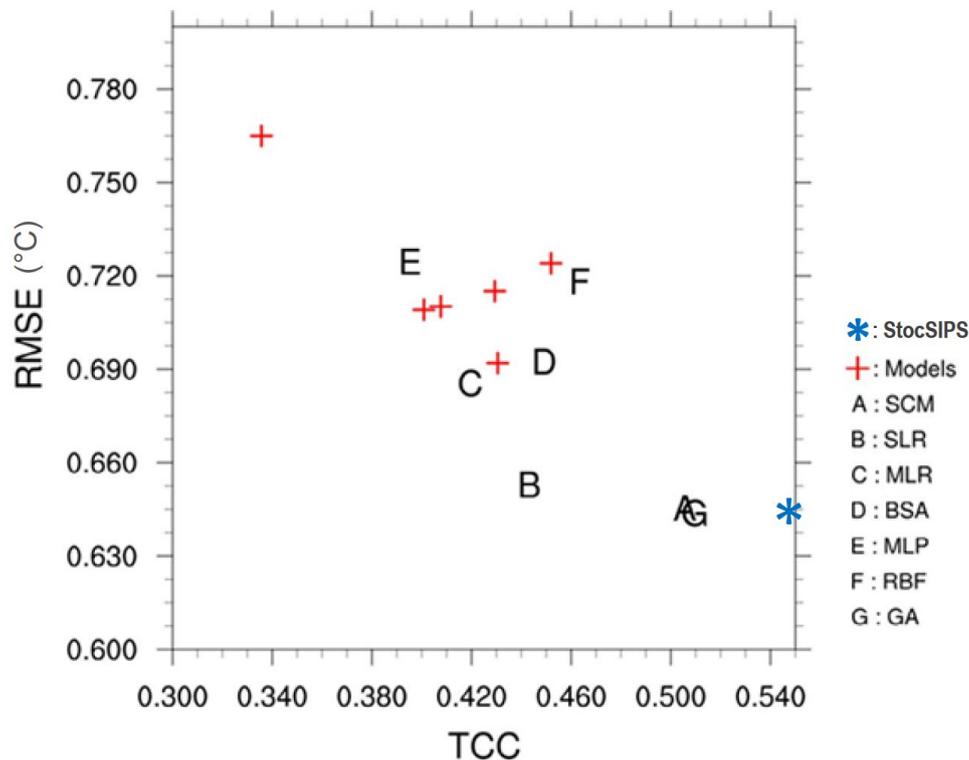


Fig. 18 4-season-averaged RMSE vs. TCC for the six individual models used in (Kim et al. 2020) (red crosses), the six MME combinations (letters) and StocSIPS (blue asterisk). This figure was adapted from Figs. 7 in (Kim et al. 2020) to include the StocSIPS scores. This is an open access article distributed under the terms of the Creative Commons CC BY license, which permits unrestricted use, distribution, and reproduction in any medium.

616 StocSIPS models the temperature series at each grid point independently as a superposition of a periodic signal corresponding
 617 to the annual cycle, a low-frequency deterministic trend from anthropogenic forcings and a high-frequency stochastic natural
 618 variability component. The annual cycle can be estimated directly from the data and is assumed constant in the future, at least
 619 for horizons of a few years. The anthropogenic component is represented as a linear response to equivalent CO₂ forcing and
 620 can be projected very accurately one year into the future by using only one parameter: the climate sensitivity, itself obtained
 621 from linear regression with historical emissions. Finally, the natural variability is modeled as a discrete-in-time fGn process
 622 which is completely determined by the variance and the fluctuation exponent. That gives a total of only three parameters for
 623 each grid point for modeling and predicting the surface temperature. Those parameters are quite stable and can be estimated
 624 with good accuracy from past data. The same procedure could be extended to any other field assuming it satisfies the
 625 Gaussianity and the scaling behaviour of the fluctuations.

626 One evident question that arises from our treatment is why not to exploit the teleconnections in the temperature field to improve
 627 the forecast instead of predicting each series independently. A detailed analysis was performed in that matter. The details are
 628 given elsewhere (Del Rio Amador and Lovejoy 2021), but the main conclusion is that, by exploiting the correlations in the
 629 temperature series, improvements on the MSSS values of only 1-2% are possible. This is in the noise level of our current

630 predictions, so in that sense, the forecast of the individual series is nearly optimal. The reason for the small improvement when
631 spatial correlations are used is that StocSIPS is a past value problem and the spatial correlations are effectively already included
632 in the past data.

633 Although we mentioned that the fGn with fluctuation exponent in the range $-1/2 < H < 0$ is a good model for the natural
634 variability, a distinction must be made for most of the tropical ocean region, for which a positive fluctuation exponent was
635 found. Instead of using the fGn model there, we must use the general fRn model or its high frequency fBm approximation with
636 $0 < H < 1$. It is significant that within this tropical ocean region, only in the more predictable region that is associated with
637 the ENSO phenomenon we obtain fluctuation exponents in the range $1/2 < H < 1$, whose fBm approximation has persistent
638 (positively correlated) increments.

639 It is surprising that by using only three parameters for each location (the fluctuation exponent, H , the standard deviation, σ_T ,
640 and the transient climate sensitivity, $\lambda_{2 \times \text{CO}_2 \text{eq}}$), we can build a model that accurately describes the temperature field. The
641 adequacy of the model was verified by testing the whiteness of the residual innovations and validating the theoretically
642 expected scores (if the model were perfect) vs. the actual hindcasts results. This also implies that, for probabilistic forecast,
643 StocSIPS is a nearly perfectly reliable system without need of recalibration of the forecast probability distribution.

644 The hindcast verification results show that the skill is generally greater over the ocean than over land, in particular over the
645 more persistent tropical ocean region. One of the implications of the scaling that was verified is that the 3-month average
646 forecast has the same skill as the one month ahead monthly prediction. This is possible because although the horizon is further
647 in the future, the seasonal forecast is for a longer (3 month) average. For scaling processes, the two effects exactly compensate.
648 The seasonal predictions show a decreased skill in the polar regions during the summer. A possible explanation for this
649 seasonality is that, when removing the annual cycle and the trend associated with the anthropogenic warming, we assumed
650 that both were statistically independent. This is not true for the polar region. While for the rest of the planet the anthropogenic
651 temperature response increases uniformly for every month following the increasing CO_2 concentrations, in the poles during
652 the summer, the temperature is tied to the freezing point of water. This spurious seasonality introduced in the preprocessing of
653 the data, can be corrected in future versions of StocSIPS to improve the global forecasts.

654 Besides this seasonality near the poles, the globally averaged skill score values are also lower during the boreal winter. This
655 can be explained by the asymmetric distribution of land mass between the northern and the southern hemispheres and the fact
656 that the atmospheric temperature near the surface is more stable over ocean than over land. Further improvements in the model
657 may be possible using recalibration of the individual forecasts for every season.

658 Although the purpose of this paper is not to make a detailed and exhaustive comparison with other long-term prediction models,
659 it is important to at least show a preliminary comparison with already published skill scores from other models to assess the
660 advantages and shortcomings of StocSIPS. The evaluation against seven different MME combination methods using six
661 models from the Lead Centers for Long-range forecasts published by (Kim et al. 2020), showed that the skill scores obtained
662 with StocSIPS are comparable (or better in the case of the temporal correlation coefficient) than the best MME combination
663 (which has larger skill than any individual ensemble member). This is in agreement with the previous results in (Del Rio

664 Amador and Lovejoy 2019) that show that StocSIPS outperformed the Canadian MME (CanSIPS) for all but the first month
665 of the forecast. This preliminary comparison for seasonal forecast validates StocSIPS as a good alternative and a
666 complementary option to conventional numerical models.

667 StocSIPS and GCMs are based on entirely different approaches. While the GCMs only take the initial state of the system (with
668 perturbations to produce multiple ensemble realizations), they do exploit all possible interactions with other atmospheric
669 variables and other locations to produce their forecast through the integration of the dynamical equations. Conversely,
670 StocSIPS neglects all the spatial (and other variables) relations to produce forecasts based on the past states at any single
671 location by exploiting the large memory of the system. Another way to view this is that for forecasts, GCMs are initial value
672 models that generate many “stochastic” realizations of the state of the atmosphere, whereas StocSIPS is effectively a “past
673 value problem” that directly estimates the most probable future state (conditional expectation).

674 Although there is no evident mechanism that explains how the distant past affects the current state of the system, model
675 reduction as explained by the Mori-Zwanzig formalism (Mori 1965; Zwanzig 1973, 2001; Gottwald et al. 2017) shows that if
676 we only look at one part of the system (e.g. the temperature at a given location), memory effects arise. All the interactions
677 coming from other degrees of freedom are embedded in the past values. Recent works (Lovejoy et al. 2015; Lovejoy 2019a,
678 2020; Lovejoy et al. 2021) hypothesize that, for the case of temperature, scaling behaviour are a result of a hierarchy of energy
679 storage mechanisms acting at different temporal and spatial scales.

680 In a recent publication (Del Rio Amador and Lovejoy 2021), StocSIPS was extended to the multivariate case (m-StocSIPS),
681 to include and realistically reproduce all the space-time cross-correlation structure. It was shown that, although large spatial
682 correlations exist in the temperature field, the optimal predictor of the temperature at a given location is obtained from its own
683 past if long enough time series are given. These cross-correlations “were already used” to build that past. This means that the
684 predictions given here (in the univariate StocSIPS version) are optimal in this stochastic framework. Nevertheless, the fact that
685 the GCMs remain “deterministic” up to approximately 1–2 years over the oceans (mostly in the tropics) and in the poles, where
686 having a dynamic sea ice model is apparently crucial for subseasonal to seasonal forecasts (Zampieri et al. 2018), suggests that
687 StocSIPS can be combined with GCM outputs to produce a single hybrid forecasting system that improves on both.

688 **Appendix A: Basic Theory for fGn Processes**

689 **i. Continuous-in-time fGn**

690 In DRAL, the stochastic natural variability component of the globally averaged temperature was represented as an fGn process.
691 The main properties of fGn relevant for the present paper are summarized in the following.

692 An fGn process at resolution τ (the scale at which the series is averaged) has the following integral representation:

$$693 \quad T_\tau(t) = \frac{1}{\tau} \frac{c_H \sigma_T}{\Gamma(H + 3/2)} \left[\int_{-\infty}^t (t-t')^{H+1/2} \gamma(t') dt' - \int_{-\infty}^{t-\tau} (t-\tau-t')^{H+1/2} \gamma(t') dt' \right], \quad (A1)$$

694 where $\gamma(t)$ is a unit Gaussian δ -correlated white noise process with $\langle \gamma(t) \rangle = 0$ and $\langle \gamma(t)\gamma(t') \rangle = \delta(t - t')$ [$\delta(x)$ is the Dirac
695 function], $\Gamma(x)$ is the Euler gamma function, σ_T is the ensemble standard deviation (for $\tau = 1$) and

$$696 \quad c_H^2 = \frac{\pi}{2 \cos(\pi H) \Gamma(-2 - 2H)}. \quad (\text{A2})$$

697 This is the canonical value for the constant c_H that was chosen to make the expression for the statistics particularly simple. In
698 particular, the variance is $\langle T_\tau(t)^2 \rangle = \sigma_T^2 \tau^{2H}$ for all t , where $\langle \cdot \rangle$ denotes ensemble (infinite realizations) averaging. The
699 parameter H , with $-1 < H < 0$, is the fluctuation exponent of the corresponding fractional Gaussian noise process, the Hurst
700 exponent, $H' = H + 1$. Fluctuation exponents are used due to their wider generality; they are well defined even for strongly
701 intermittent non-Gaussian multifractal processes and they can be any real value. For a discussion, see page 643 in (Lovejoy et
702 al. 2015).

703 Equation (A1) can be interpreted as the smoothing of the fractional integral of a white noise process or as the power-law
704 weighted average of past innovations, $\gamma(t)$. This power-law weighting accounts for the memory effects in the temperature
705 series. The closer the fluctuation exponent is to zero, the larger is the influence of past values on the current temperature. This
706 is evidenced by the behaviour of the autocorrelation function:

$$707 \quad R_H(\Delta t) = \frac{\langle T_\tau(t) T_\tau(t + \Delta t) \rangle}{\langle T_\tau(t)^2 \rangle} = \frac{1}{2} \left(\left| \frac{\Delta t}{\tau} + 1 \right|^{2H+2} + \left| \frac{\Delta t}{\tau} - 1 \right|^{2H+2} - 2 \left| \frac{\Delta t}{\tau} \right|^{2H+2} \right), \quad (\text{A3})$$

708 for $|\Delta t| \geq \tau$. In particular, for $\Delta t \gg \tau$ we obtain:

$$709 \quad R_H(\Delta t) \approx (H + 1)(2H + 1) \left(\frac{\Delta t}{\tau} \right)^{2H}, \quad (\text{A4})$$

710 which has a power-law behaviour with the same exponent as the average squared fluctuation and due to the Wiener–Khinchin
711 theorem, it implies the spectrum exponent $\beta = 1 + 2H$. For more details on fGn processes see (Mandelbrot and Van Ness
712 1968; Gripenberg and Norros 1996; Biagini et al. 2008).

713 **ii. Discrete-in-time fGn**

714 A detailed explanation of the theory for modeling and predicting using the discrete version of fGn processes was presented in
715 DRAL; the main results are summarized next. The analogue of Eq. (A1) in the discrete case for a finite series, $\{T_t\}_{t=1, \dots, N}$, with
716 length N and zero mean is:

$$717 \quad T_t = \sum_{j=1}^t m_j \gamma_{t+1-j} = m_{t1} \gamma_t + \dots + m_{tt} \gamma_1, \quad (\text{A5})$$

718 for $t = 1, \dots, N$, where $\{\gamma_t\}_{t=1, \dots, N}$ is a discrete white noise process and the coefficients m_{ij} are the elements of the lower
719 triangular matrix $\mathbf{M}_{H, \sigma_T}^N$ given by the Cholesky decomposition of the autocovariance matrix, $\mathbf{C}_{H, \sigma_T}^N = \sigma_T^2 [R_H(i - j)]_{i, j=1, \dots, N}$:

$$720 \quad \mathbf{C}_{H, \sigma_T}^N = \mathbf{M}_{H, \sigma_T}^N \left(\mathbf{M}_{H, \sigma_T}^N \right)^T, \quad (\text{A6})$$

721 with $m_{ij} = 0$ for $j > i$ (we assume $\tau = 1$ is the smallest scale in our system). The superscript T denotes transpose operation.

722 In vector form, Eq. (A5) can be written as:

$$723 \quad \mathbf{T}_N = \mathbf{M}_{H, \sigma_T}^N \boldsymbol{\gamma}_N \quad (\text{A7})$$

724 Equations (A5-A7) can be used to create synthetic samples of fGn with a given length N , autocorrelation function given by

725 Eq. (A3) and set of parameters $\sigma_T > 0$ and $-1 < H < 0$ (the mean of the series is always assumed equal to zero). Conversely,

726 given an actual temperature series with vector $\mathbf{T}_N = [T_1, \dots, T_N]^T$, we can estimate the parameters σ_T and H using the

727 maximum likelihood method (details are given in Appendix 1 of DRAL) and we can verify that it could be well approximated

728 by an fGn model by inverting Eq. (A7) and obtaining the residual vector of innovations:

$$729 \quad \boldsymbol{\gamma}_N = \left(\mathbf{M}_{H, \sigma_T}^N \right)^{-1} \mathbf{T}_N. \quad (\text{A8})$$

730 If the model provides a good description of the data, the residual vector $\boldsymbol{\gamma}_N = [\gamma_1, \dots, \gamma_N]^T$ is a white noise, i.e. the elements

731 should be NID(0,1) with autocorrelation function $\langle \gamma_i \gamma_j \rangle = \delta_{ij}$ (δ_{ij} is the Kronecker delta and NID(0,1) stands for Normally

732 and Independently Distributed with mean 0 and variance 1). It is worth mentioning that a white noise process is a particular

733 case of fGn with $H = -1/2$.

734 **iii. fRn Correlation Function for $0 < H < 1$**

735 The fractional Relaxation noise (fRn) process was introduced in (Lovejoy 2019b) generalizing both fGn, fBm and Ornstein-

736 Uhlenbeck processes. For short time scales (compared to some characteristic relaxation time, τ_r) and for exponents $-1/2 <$

737 $H < 0$, the fRn is close to an fGn process. For fluctuation exponents in the range $0 < H < 1$ the high-frequency approximation

738 to fRn is no longer an fGn process. In this case, to leading order, the correlation function is:

$$739 \quad R_{\text{fRn}}(\Delta t) = 1 - A_H \left(\frac{\Delta t}{\tau_r} \right)^{2H} + O \left(\frac{\Delta t}{\tau_r} \right)^{3H+1/2}; \quad \Delta t < \tau_r; \quad 0 < H < 1 \quad (\text{A9})$$

740 where τ_r is the relaxation time and A_H is an H -dependent numerical factor [see (Lovejoy 2019b)]. The same correlation

741 function was obtained by (Delignières 2015) as an approximation to short segments of discrete-in-time fractional Brownian

742 motion (fBm) process that is the integral of an fGn process (but with H increased by 1). This shows that although fBm is

743 nonstationary, short segments approximate (the stationary) fRn process. When $0 < H < 1$, fBm is a high-frequency

744 approximation to an fRn process.

745 **iv. Prediction**

746 In DRAL it was shown that, if $\{T_t\}_{t < 0}$ is an fGn process, the optimal k -steps predictor for T_k ($k > 0$), based on a finite number,

747 m (memory), of past values, is given by:

$$748 \quad \hat{T}_k = \sum_{j=-m}^0 \phi_j(k) T_j = \phi_{-m}(k) T_{-m} + \dots + \phi_0(k) T_0, \quad (\text{A10})$$

749 where the vector of predictor coefficients, $\boldsymbol{\phi}(k) = [\phi_{-m}(k), \dots, \phi_0(k)]^T$, satisfies the Yule-Walker equations:

750
$$\mathbf{R}_H \boldsymbol{\phi}(k) = \mathbf{r}_H(k), \quad (\text{A11})$$

751 with the vector $\mathbf{r}_H(k) = [R_H(k-i)]_{i=-m, \dots, 0}^T = [R_H(m+k), \dots, R_H(k)]^T$ and $\mathbf{R}_H = [R_H(i-j)]_{i,j=-t, \dots, 0}$ being the
 752 autocorrelation matrix (see Eq. (A3)). In those regions with consecutive values positively correlated (blue regions in Fig. 4a
 753 with $-1/2 < H < 0$ or the increments in the yellow region with $1/2 < H < 1$), the elements $R_H(\Delta t)$ are obtained from Eq.
 754 (A3). In the places with consecutive increments negatively correlated, where $0 < H < 1/2$ (red in Fig. 4a), instead of
 755 forecasting the fGn increments, we forecast directly the fRn process and we get the elements $R_H(\Delta t)$ from Eq. (A9). To use
 756 this autocorrelation for fRn, we estimate the constant A_H in Eq. (A9) for each location by fitting the empirical autocorrelation
 757 function.

758 The root mean square error (RMSE) for the predictor at a future time k , using a memory of m values, is defined as:

759
$$\text{RMSE}(k, m) = \sqrt{\left\langle \left[T_k - \hat{T}_k(m) \right]^2 \right\rangle}. \quad (\text{A12})$$

760 Following the results presented in DRAL and using that, for positive H the fRn is the integral of the corresponding fGn process,
 761 we obtain the following analytical expression for the RMSE of the predictor of the natural variability component:

762
$$\text{RMSE}_{\text{nat}}^{\text{theory}}(k) = \text{RMSE}(k, m, \sigma_T, H) = \begin{cases} \sigma_T \sqrt{1 - \mathbf{r}_H(k)^T (\mathbf{R}_H)^{-1} \mathbf{r}_H(k)} ; & \text{for } -1/2 < H < 0 \\ \sigma_T k^H \sqrt{1 - \mathbf{r}_{H-1}(1)^T (\mathbf{R}_{H-1})^{-1} \mathbf{r}_{H-1}(1)} ; & \text{for } 0 < H < 1 \end{cases}. \quad (\text{A13})$$

763 For a given forecast horizon, k , the RMSE only depends on the parameters σ_T and H , and the memory used, m . In Fig. 3 of
 764 DRAL it was shown that only a few past datapoints are needed as memory to obtain an error approaching – with more than
 765 95% agreement – the asymptotical value corresponding to $m = \infty$, for all possible values of H .

766 The theoretical mean square skill score (MSSS), is defined as:

767
$$\text{MSSS}(k) = 1 - \frac{\left\langle \left[T(k) - \hat{T}(k) \right]^2 \right\rangle}{\left\langle T(k)^2 \right\rangle} \quad (\text{A14})$$

768 (the reference forecast is the mean of the series, assumed equal to zero here).

769 From the definition of the RMSE, Eq. (A12), we obtain the theoretical value for fGn:

770
$$\text{MSSS}_{\text{nat}}^{\text{theory}}(k) = \text{MSSS}(k, m, H) = 1 - \frac{\text{RMSE}(k, m, \sigma_T, H)^2}{\sigma_T^2} \quad (\text{A15})$$

771 or, replacing Eq. (A13) for $-1/2 < H < 0$:

772
$$\text{MSSS}(k, m, H) = \mathbf{r}_H(k)^T (\mathbf{R}_H)^{-1} \mathbf{r}_H(k) = \boldsymbol{\phi}(k) \cdot \mathbf{r}_H(k). \quad (\text{A16})$$

773 In Fig. 19 we show graphs of the theoretical MSSS as a function of H for different values of k . A memory $m = 50$ was used
 774 for computing the MSSS. As expected, the skill decreases as the forecast horizon increases. For $H = -0.5$, the fGn process is
 775 a white noise process and $\text{MSSS} = 0$. The skill increases with H and (with infinite past data) the process becomes perfectly
 776 predictable when $H \rightarrow 0$.

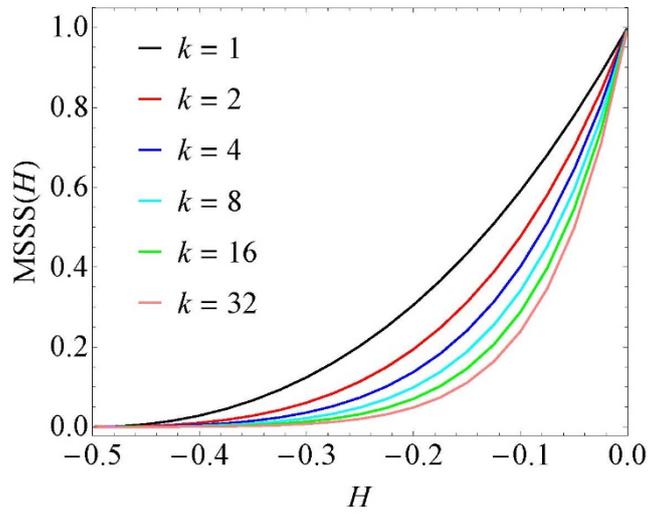


Fig. 19 Graphs of the theoretical MSSS (Eq. (A16)) as a function of H for different values of k . A memory $m = 50$ was used for computing the MSSS.

777 Appendix B: Verification Metrics

778 i. Definitions

779 The verification metrics used in this paper were defined following the recommendations in the Standardized verification system
 780 for long-range forecasts (SVS-LRF) for the practical details of producing and exchanging appropriate verification scores
 781 (WMO 2010b, a). Let $x_i(t)$ and $f_i(t)$, ($t = 1, \dots, N$) denote time series of observations and forecasts, respectively, for a grid
 782 point i over the period of verification (POV) with N time steps. Then, their averages for the POV, \bar{x}_i and \bar{f}_i and their sample
 783 variances $s_{x_i}^2$ and $s_{f_i}^2$ are given by:

$$\begin{aligned}
 \bar{x}_i &= \frac{1}{N} \sum_{t=1}^N x_i(t), & \bar{f}_i &= \frac{1}{N} \sum_{t=1}^N f_i(t) \\
 s_{x_i}^2 &= \frac{1}{N} \sum_{t=1}^N [x_i(t) - \bar{x}_i]^2, & s_{f_i}^2 &= \frac{1}{N} \sum_{t=1}^N [f_i(t) - \bar{f}_i]^2.
 \end{aligned}
 \tag{B1}$$

785 The mean square error (MSE) of the forecast for grid point i is:

$$\text{MSE}_i = \frac{1}{N} \sum_{t=1}^N [f_i(t) - x_i(t)]^2
 \tag{B2}$$

787 and the root mean square error (RMSE) is:

$$\text{RMSE}_i = \sqrt{\text{MSE}_i}.
 \tag{B3}$$

789 For leave-one-out cross-validated data in the POV (WMO 2010a), the MSE of climatology forecasts is:

$$\text{MSE}_{C_i} = \left(\frac{N}{N-1} \right)^2 s_{x_i}^2.
 \tag{B4}$$

791 The mean square skill score (MSSS) for grid point i , taking as reference the climatology forecast, is defined as:

$$792 \quad \text{MSSS}_i = 1 - \frac{\text{MSE}_i}{\text{MSE}_{C_i}}. \quad (\text{B5})$$

793 The temporal correlation coefficient (TCC) is:

$$794 \quad \text{TCC}_i = \frac{\frac{1}{N} \sum_{t=1}^N [x_i(t) - \bar{x}_i][f_i(t) - \bar{f}_i]}{s_{x_i} s_{f_i}}. \quad (\text{B6})$$

795 Both the MSE_i and the TCC_i are computed using temporal averages for a given location i , conversely, the anomaly pattern
796 correlation coefficient (ACC) (Jolliffe and Stephenson 2011) is defined using spatial averages for a given time t :

$$797 \quad \text{ACC}(t) = \frac{\sum_{i=1}^n \cos \theta_i [x'_i(t) - \langle x'(t) \rangle][f'_i(t) - \langle f'(t) \rangle]}{\sqrt{\sum_{i=1}^n \cos \theta_i [x'_i(t) - \langle x'(t) \rangle]^2} \sqrt{\sum_{i=1}^n \cos \theta_i [f'_i(t) - \langle f'(t) \rangle]^2}}, \quad (\text{B7})$$

798 where n is the number of grid points, θ_i is the latitude at location i , $x'_i(t)$ and $f'_i(t)$ are observation and forecast anomalies for
799 the POV, respectively, and the spatial averages $\langle x'(t) \rangle$ and $\langle f'(t) \rangle$ are given by:

$$800 \quad \langle x'(t) \rangle = \frac{\sum_{i=1}^n \cos \theta_i x'_i(t)}{\sum_{i=1}^n \cos \theta_i}, \quad \langle f'(t) \rangle = \frac{\sum_{i=1}^n \cos \theta_i f'_i(t)}{\sum_{i=1}^n \cos \theta_i} \quad (\text{B8})$$

801 **ii. Averaged scores**

802 To take the average of nonlinear scores, they should be transformed so the corresponding variables are Gaussian. The spatial
803 average RMSE (considering the area factor) is:

$$804 \quad \langle \text{RMSE} \rangle = \sqrt{\frac{\sum_{i=1}^n \text{MSE}_i \cos \theta_i}{\sum_{i=1}^n \cos \theta_i}}. \quad (\text{B9})$$

805 Similarly, the average MSSS is:

$$806 \quad \langle \text{MSSS} \rangle = 1 - \frac{\sum_{i=1}^n \text{MSE}_i \cos \theta_i}{\sum_{i=1}^n \text{MSE}_{C_i} \cos \theta_i}. \quad (\text{B10})$$

807 For the correlation coefficients, the Fisher Z-transform must be taken first. This is defined as:

$$808 \quad Z(r) = \frac{1}{2} \ln \left(\frac{1+r}{1-r} \right) = \tanh^{-1} r \quad (\text{B11})$$

809 The spatial average TCC is the defined as:

$$810 \quad \langle \text{TCC} \rangle = Z^{-1} \left[\frac{\sum_{i=1}^n Z(\text{TCC}_i) \cos \theta_i}{\sum_{i=1}^n \cos \theta_i} \right] \quad (\text{B12})$$

811 and the temporal average ACC is

$$812 \quad \langle \text{ACC} \rangle = Z^{-1} \left\{ \frac{1}{N} \sum_{t=1}^N Z[\text{ACC}(t)] \right\}. \quad (\text{B13})$$

813 **iii. Orthogonality principle and MSSS decomposition**

814 The MSSS (Eq. (B5)), can be expanded for leave-one-out cross-validated forecasts (Murphy 1988). Using Eqs. (B1), (B2),
815 (B4) and (B6) in (B5), we obtain:

$$816 \quad \text{MSSS}_i = \left\{ 2 \frac{s_{\hat{f}_i}}{s_{x_i}} \text{TCC}_i - \left(\frac{s_{\hat{f}_i}}{s_{x_i}} \right)^2 - \left(\frac{[\bar{f}_i - \bar{x}_i]}{s_{x_i}} \right)^2 + \frac{2N-1}{(N-1)^2} \right\} / \left\{ 1 + \frac{2N-1}{(N-1)^2} \right\}. \quad (\text{B14})$$

817 This equation gives a relation between the MSSS and the TCC. For forecasts with the same variance as that of observations and
818 no overall bias, the MSSS is only positive (MSE lower than for climatology) if the TCC is larger than approximately 0.5.

819 A more simplified relation can be obtained in our case for the prediction of the detrended anomalies (natural variability). As
820 we mentioned in Appendix Aiv, the predictor (Eq. (A10)) is built in such a way that the coefficients satisfy the Yule Walker
821 equations, which are derived from the orthogonality principle (Wold 1938; Brockwell and Davis 1991; Hipel and McLeod
822 1994; Palma 2007; Box et al. 2008). This principle states that the error of the optimal predictor, $e_i(t) = x_i(t) - f_i(t)$ (in a
823 mean square error sense) is orthogonal to any possible estimator:

$$824 \quad \langle e_i(t) f_i(t) \rangle = 0. \quad (\text{B15})$$

825 From this ensemble average condition, we get the analytical expressions for the coefficients as a function of the fluctuation
826 exponent, H , for the fGn process. If the model realistically describes the actual temperature anomalies, then the condition Eq.
827 (B15) can be approximated by the temporal average in the POV:

$$828 \quad \frac{1}{N} \sum_{t=1}^N [e_i(t) f_i(t)] = 0. \quad (\text{B16})$$

829 or

$$830 \quad \frac{1}{N} \sum_{t=1}^N [x_i(t) - f_i(t)] f_i(t) = 0. \quad (\text{B17})$$

831 from which:

$$832 \quad \frac{1}{N} \sum_{t=1}^N [x_i(t) f_i(t)] = \frac{1}{N} \sum_{t=1}^N f_i(t)^2. \quad (\text{B18})$$

833 For $\bar{x}_i = \bar{f}_i = 0$, dividing by the product $s_{x_i}s_{f_i}$ and using Eqs. (B1) and (B6), we can rewrite Eq. (B18) as:

$$834 \quad TCC_i = \frac{s_{f_i}}{s_{x_i}}. \quad (B19)$$

835 Using this ratio in Eq. (B14) we finally obtain:

$$836 \quad MSSS_i = \frac{TCC_i^2 + \frac{2N-1}{(N-1)^2}}{1 + \frac{2N-1}{(N-1)^2}}. \quad (B20)$$

837 A more detailed analysis gives the same expression with the weaker condition of overall unbiased estimates $\bar{x}_i - \bar{f}_i = 0$ (not
838 necessarily each of them must be zero).

839 In our case, for the forecast of the detrended anomalies (natural variability) at monthly resolution in the POV 1951-2019 ($N =$
840 828 months), the N -dependent term in Eq. (B20) is negligible:

$$841 \quad \frac{2N-1}{(N-1)^2} \approx 0.0024, \quad (B21)$$

842 so, with good approximation we obtain:

$$843 \quad MSSS_i \approx TCC_i^2. \quad (B22)$$

844 The orthogonality principle, Eq. (B17) (or equivalently, Eq. (B19) or Eq. (B22)), is the condition that maximizes the MSSS. In
845 our case, where the autoregressive coefficients in our predictor are analytical functions of only one parameter (H), if Eq. (B22)
846 is verified then our predictor is optimal in a mean square error sense and our model is suitable for describing the natural
847 temperature variability.

848 **References**

- 849 Biagini F, Hu Y, Øksendal B, Zhang T (2008) Stochastic Calculus for Fractional Brownian Motion and Applications. Springer
850 London, London
- 851 Blender R, Fraedrich K, Hunt B (2006) Millennial climate variability: GCM-simulation and Greenland ice cores. Geophys Res
852 Lett 33:L04710. doi: 10.1029/2005GL024919
- 853 Box GEP, Jenkins GM, Reinsel GC (2008) Time Series Analysis. Wiley
- 854 Brockwell PJ, Davis RA (1991) Time Series: Theory and Methods. Springer New York, New York, NY
- 855 Christensen HM, Berner J, Coleman DRB, Palmer TN (2017) Stochastic Parameterization and El Niño–Southern Oscillation.
856 J Clim 30:17–38. doi: 10.1175/JCLI-D-16-0122.1
- 857 Christensen HM, Moroz IM, Palmer TN (2015) Evaluation of ensemble forecast uncertainty using a new proper score:
858 Application to medium-range and seasonal forecasts. Q J R Meteorol Soc 141:538–549. doi: 10.1002/qj.2375

859 Clarke DC, Richardson M (2020) The Benefits of Continuous Local Regression for Quantifying Global Warming. *Earth Sp*
860 *Sci Open Arch*. doi: 10.1002/essoar.10502294.1

861 Cleveland WS, Devlin SJ (1988) Locally Weighted Regression: An Approach to Regression Analysis by Local Fitting. *J Am*
862 *Stat Assoc* 83:596. doi: 10.2307/2289282

863 Davini P, von Hardenberg J, Corti S, et al (2017) Climate SPHINX: evaluating the impact of resolution and stochastic physics
864 parameterisations in the EC-Earth global climate model. *Geosci Model Dev* 10:1383–1402. doi: 10.5194/gmd-10-1383-
865 2017

866 Del Rio Amador L, Lovejoy S (2019) Predicting the global temperature with the Stochastic Seasonal to Interannual Prediction
867 System (StocSIPS). *Clim Dyn* 53:4373–4411. doi: 10.1007/s00382-019-04791-4

868 Del Rio Amador L, Lovejoy S (2021) Long-range Forecasting as a Past Value Problem: Untangling Correlations and Causality
869 with Scaling. *Geophys Res Lett* in review: doi: 10.1002/essoar.10505160.1

870 Delignières D (2015) Correlation Properties of (Discrete) Fractional Gaussian Noise and Fractional Brownian Motion. *Math*
871 *Probl Eng* 2015:1–7. doi: 10.1155/2015/485623

872 Fisher RA (1915) Frequency Distribution of the Values of the Correlation Coefficient in Samples from an Indefinitely Large
873 Population. *Biometrika* 10:507. doi: 10.2307/2331838

874 Franzke C (2012) Nonlinear Trends, Long-Range Dependence, and Climate Noise Properties of Surface Temperature. *J Clim*
875 25:4172–4183. doi: 10.1175/JCLI-D-11-00293.1

876 Franzke CLE, O’Kane TJ, Berner J, et al (2015) Stochastic climate theory and modeling. *Wiley Interdiscip Rev Clim Chang*
877 6:63–78. doi: 10.1002/wcc.318

878 Gneiting T, Raftery AE, Westveld AH, Goldman T (2005) Calibrated Probabilistic Forecasting Using Ensemble Model Output
879 Statistics and Minimum CRPS Estimation. *Mon Weather Rev* 133:1098–1118. doi: 10.1175/mwr2904.1

880 Gottwald GA, Crommelin DT, Franzke CLE (2017) Stochastic Climate Theory. In: Franzke CLE, OKane TJ (eds) *Nonlinear*
881 *and Stochastic Climate Dynamics*. Cambridge University Press, Cambridge, pp 209–240

882 Graham R, Yun W, Kim J, et al (2011) Long-range forecasting and the Global Framework for Climate Services. *Clim Res*
883 47:47–55. doi: 10.3354/cr00963

884 Gripenberg G, Norros I (1996) On the prediction of fractional Brownian motion. *J Appl Probab* 33:400–410. doi:
885 10.1017/S0021900200099812

886 Hasselmann K (1976) Stochastic climate models Part I. Theory. *Tellus* 28:473–485. doi: 10.1111/j.2153-3490.1976.tb00696.x

887 Hébert R, Lovejoy S (2018) Regional Climate Sensitivity- and Historical-Based Projections to 2100. *Geophys Res Lett*
888 45:4248–4254. doi: 10.1002/2017GL076649

889 Hersbach H (2000) Decomposition of the Continuous Ranked Probability Score for Ensemble Prediction Systems. *Weather*
890 *Forecast* 15:559–570. doi: 10.1175/1520-0434(2000)015<0559:DOTCRP>2.0.CO;2

891 Hipel KW, McLeod AI (1994) Time Series Modelling of Water Resources and Environmental Systems. In: Hipel KW, McLeod
892 AI (eds) *Time Series Modelling of Water Resources and Environmental Systems*. Elsevier, pp iii–xxxvii, 1–1013

893 Jolliffe IT, Stephenson DB (2011) *Forecast Verification: A Practitioner's Guide in Atmospheric Science*, 2nd edn. John Wiley
894 and Sons

895 Kalnay E, Kanamitsu M, Kistler R, et al (1996) The NCEP/NCAR 40-Year Reanalysis Project. *Bull Am Meteorol Soc* 77:437–
896 471. doi: 10.1175/1520-0477(1996)077<0437:TNYRP>2.0.CO;2

897 Keller JD, Hense A (2011) A new non-Gaussian evaluation method for ensemble forecasts based on analysis rank histograms.
898 *Meteorol Zeitschrift* 20:107–117. doi: 10.1127/0941-2948/2011/0217

899 Kim G, Ahn J, Kryjov VN, et al (2020) Assessment of MME methods for seasonal prediction using WMO LC-LRFMME
900 hindcast dataset. *Int J Climatol* *joc.6858*. doi: 10.1002/joc.6858

901 Koscielny-Bunde E, Bunde A, Havlin S, et al (1998) Indication of a Universal Persistence Law Governing Atmospheric
902 Variability. *Phys Rev Lett* 81:729–732. doi: 10.1103/PhysRevLett.81.729

903 Kryjov VN, Kang H-W, Nohara D, et al (2006) Assessment of the climate forecasts produced by individual models and MME
904 methods. APCC Technical Report 2006, APEC Climate Center. Busan, South Korea

905 Leutbecher M, Palmer TN (2008) Ensemble forecasting. *J Comput Phys* 227:3515–3539. doi: 10.1016/j.jcp.2007.02.014

906 Lovejoy S (2014) Scaling fluctuation analysis and statistical hypothesis testing of anthropogenic warming. *Clim Dyn* 42:2339–
907 2351. doi: 10.1007/s00382-014-2128-2

908 Lovejoy S (2020) The fractional heat equation. *Geophys Res Lett* (in review)

909 Lovejoy S (2019a) Fractional relaxation noises, motions and the fractional energy balance equation. *Nonlin Process Geophys*
910 Discuss [preprint] in review: doi: <https://doi.org/10.5194/npg-2019-39>

911 Lovejoy S (2021a) The Half-order Energy Balance Equation, Part 1: The homogeneous HEBE and long memories. *Earth Syst*
912 *Dynam* (in press): doi: <https://doi.org/10.5194/esd-2020-12>

913 Lovejoy S (2021b) The Half-order Energy Balance Equation, Part 2: The inhomogeneous HEBE and 2D energy balance
914 models. *Earth Syst Dynam Discuss* (in press): doi: <https://doi.org/10.5194/esd-2020-13>

915 Lovejoy S (2018) Spectra, intermittency, and extremes of weather, macroweather and climate. *Sci Rep* 8:12697. doi:
916 10.1038/s41598-018-30829-4

917 Lovejoy S (2019b) Fractional relaxation noises, motions and the fractional energy balance equation. *Nonlin Process Geophys*
918 Discuss in review: doi: <https://doi.org/10.5194/npg-2019-39>

919 Lovejoy S, del Rio Amador L, Hébert R (2015) The ScaLIng Macroweather Model (SLIMM): using scaling to forecast global-
920 scale macroweather from months to decades. *Earth Syst Dyn* 6:637–658. doi: 10.5194/esd-6-637-2015

921 Lovejoy S, Procyk R, Hébert R, Del Rio Amador L (2021) The Fractional Energy Balance Equation. *Q J R Meteorol Soc*
922 *qj.4005*. doi: 10.1002/qj.4005

923 Lovejoy S, Schertzer D (2010) Towards a new synthesis for atmospheric dynamics: Space–time cascades. *Atmos Res* 96:1–
924 52. doi: 10.1016/j.atmosres.2010.01.004

925 Lovejoy S, Schertzer D (2013) *The Weather and Climate: Emergent Laws and Multifractal Cascades*. Cambridge University
926 Press, Cambridge

927 Lovejoy S, Schertzer D (1986) Scale invariance in climatological temperatures and the spectral plateau. *Ann Geophys* 4B:401–
928 410

929 Lovejoy S, Schertzer D (2012a) Haar wavelets, fluctuations and structure functions: convenient choices for geophysics.
930 *Nonlinear Process Geophys* 19:513–527. doi: 10.5194/npg-19-513-2012

931 Lovejoy S, Schertzer D (2012b) Low-Frequency Weather and the Emergence of the Climate. pp 231–254

932 Mandelbrot BB, Van Ness JW (1968) Fractional Brownian Motions, Fractional Noises and Applications. *SIAM Rev* 10:422–
933 437. doi: 10.1137/1010093

934 Meinshausen M, Smith SJ, Calvin K, et al (2011) The RCP greenhouse gas concentrations and their extensions from 1765 to
935 2300. *Clim Change* 109:213–241. doi: 10.1007/s10584-011-0156-z

936 Mori H (1965) Transport, Collective Motion, and Brownian Motion. *Prog Theor Phys* 33:423–455. doi: 10.1143/PTP.33.423

937 Murphy AH (1988) Skill Scores Based on the Mean Square Error and Their Relationships to the Correlation Coefficient. *Mon*
938 *Weather Rev* 116:2417–2424. doi: 10.1175/1520-0493(1988)116<2417:SSBOTM>2.0.CO;2

939 NCEP/NCAR (2020) NCEP/NCAR Reanalysis 1. <https://psl.noaa.gov/data/gridded/data.ncep.reanalysis.html>. Accessed 3 Jan
940 2020

941 Newman M (2013) An Empirical Benchmark for Decadal Forecasts of Global Surface Temperature Anomalies. *J Clim*
942 26:5260–5269. doi: 10.1175/JCLI-D-12-00590.1

943 Newman M, Sardeshmukh PD, Winkler CR, Whitaker JS (2003) A Study of Subseasonal Predictability. *Mon Weather Rev*
944 131:1715–1732. doi: <https://doi.org/10.1175//2558.1>

945 Palma W (2007) Long-Memory Time Series. John Wiley & Sons, Inc., Hoboken, NJ, USA

946 Palmer T, Buizza R, Hagedorn R, et al (2006) Ensemble prediction: a pedagogical perspective. *ECMWF Newsl* 106:10–17.
947 doi: 10.21957/ab129056ew

948 Palmer TN (2019) Stochastic weather and climate models. *Nat Rev Phys* 1:463–471. doi: 10.1038/s42254-019-0062-2

949 Pasternack A, Bhend J, Liniger MA, et al (2018) Parametric decadal climate forecast recalibration (DeFoReSt 1.0). *Geosci*
950 *Model Dev* 11:351–368. doi: 10.5194/gmd-11-351-2018

951 Penland C, Matrosova L (1994) A Balance Condition for Stochastic Numerical Models with Application to the El Niño-
952 Southern Oscillation. *J Clim* 7:1352–1372. doi: 10.1175/1520-0442(1994)007<1352:ABCFSN>2.0.CO;2

953 Penland C, Sardeshmukh PD (1995) The optimal growth of tropical sea surface temperature anomalies. *J Clim* 8:1999–2024.
954 doi: 10.1175/1520-0442(1995)008<1999:TOGOTS>2.0.CO;2

955 Procyk R, Lovejoy S, Hébert R (2020) The Fractional Energy Balance Equation for Climate projections through 2100. *Earth*
956 *Syst Dyn* (submitted)

957 Rackow T, Juricke S (2020) Flow-dependent stochastic coupling for climate models with high ocean-to-atmosphere resolution
958 ratio. *Q J R Meteorol Soc* 146:284–300. doi: 10.1002/qj.3674

959 Rypdal K, Østvand L, Rypdal M (2013) Long-range memory in Earth’s surface temperature on time scales from months to
960 centuries. *J Geophys Res Atmos* 118:7046–7062. doi: 10.1002/jgrd.50399

961 Sardeshmukh PD, Sura P (2009) Reconciling non-Gaussian climate statistics with linear dynamics. *J Clim* 22:1193–1207. doi:
962 10.1175/2008JCLI2358.1

963 Trenberth KE (1997) The Definition of El Niño. *Bull Am Meteorol Soc* 78:2771–2777. doi: 10.1175/1520-
964 0477(1997)078<2771:TDOENO>2.0.CO;2

965 Varotsos CA, Efsthathiou MN, Cracknell AP (2013) On the scaling effect in global surface air temperature anomalies. *Atmos*
966 *Chem Phys* 13:5243–5253. doi: 10.5194/acp-13-5243-2013

967 Williams PD (2012) Climatic impacts of stochastic fluctuations in air-sea fluxes. *Geophys Res Lett* 39:. doi:
968 10.1029/2012GL051813

969 Winkler CR, Newman M, Sardeshmukh PD (2001) A linear model of wintertime low-frequency variability. Part I: Formulation
970 and forecast skill. *J Clim* 14:4474–4494. doi: 10.1175/1520-0442(2001)014<4474:ALMOWL>2.0.CO;2

971 WMO (2010a) Standardised verification system (SVS) for long-range forecasts (LRF). New attachment II-8 to the manual on
972 the GDPS. WMO-No. 485, Volume 1. Geneva, Switzerland.

973 WMO (2010b) Manual on the Global Data-processing and Forecasting System Volume I. (WMO-No. 485). Geneva,
974 Switzerland.

975 Wold H (1938) A Study in Analysis of Stationary Time Series. *J R Stat Soc*

976 Wu Z, Huang NE (2009) Ensemble Empirical Mode Decomposition: A Noise-Assisted Data Analysis Method. *Adv Adapt*
977 *Data Anal* 01:1–41. doi: 10.1142/S1793536909000047

978 Yuan N, Fu Z, Liu S (2015) Extracting climate memory using Fractional Integrated Statistical Model: A new perspective on
979 climate prediction. *Sci Rep* 4:6577. doi: 10.1038/srep06577

980 Zampieri L, Goessling HF, Jung T (2018) Bright Prospects for Arctic Sea Ice Prediction on Subseasonal Time Scales. *Geophys*
981 *Res Lett* 45:9731–9738. doi: 10.1029/2018GL079394

982 Zeiler A, Faltermeier R, Keck IR, et al (2010) Empirical Mode Decomposition - an introduction. In: *The 2010 International*
983 *Joint Conference on Neural Networks (IJCNN)*. IEEE, pp 1–8

984 Zwanzig R (2001) *Nonequilibrium Statistical Mechanics*, 1st edn. Oxford University Press, U.S.A

985 Zwanzig R (1973) Nonlinear generalized Langevin equations. *J Stat Phys* 9:215–220. doi: 10.1007/BF01008729

986

Figures

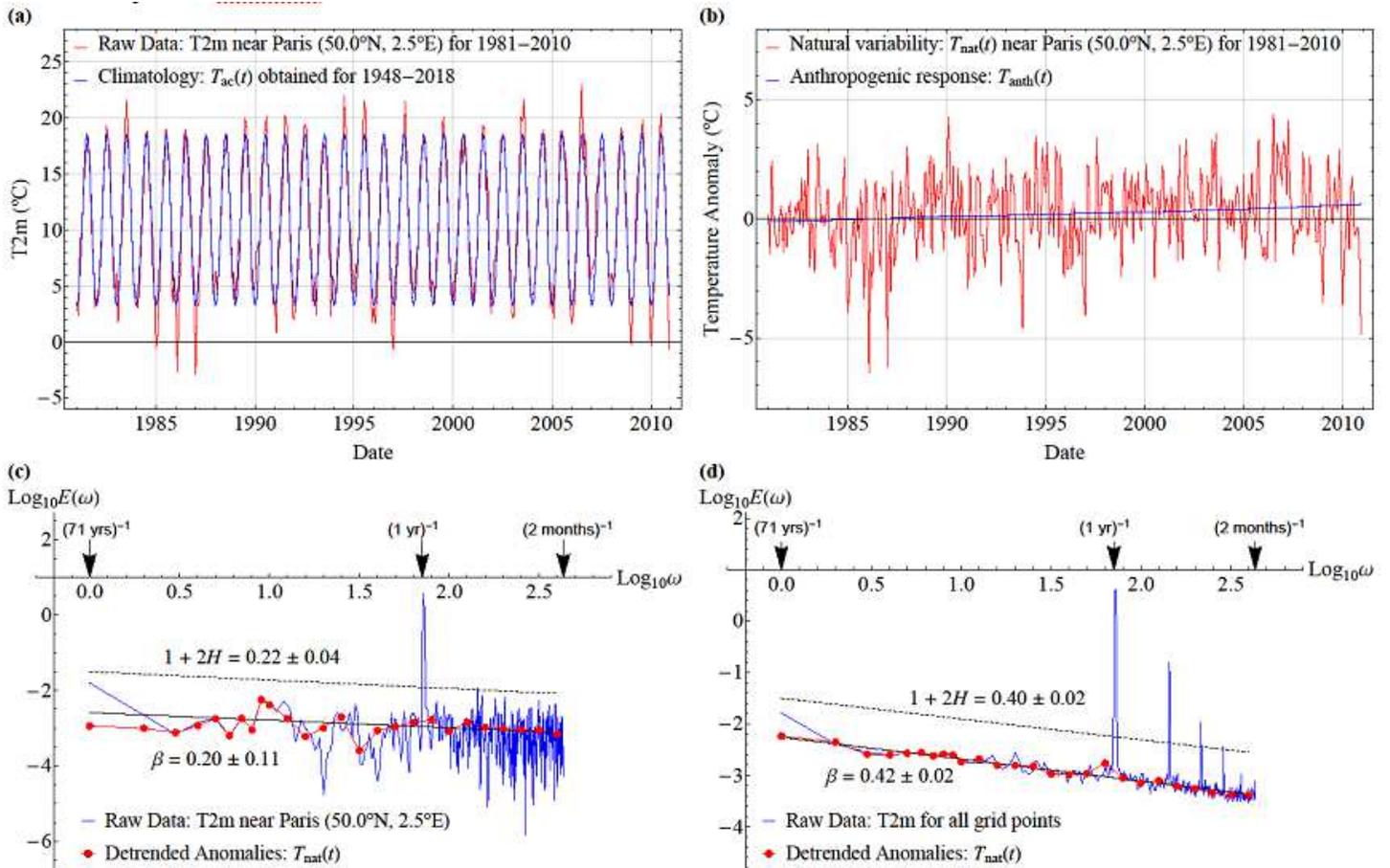


Figure 1

Example of signal pre-processing and spectra for the grid point with coordinates 50.0°N, 2.5°E (near Paris, France). (a) Raw temperature data, T (in red), and the periodic signal, T_{ac} (in blue). Only the period 1981–2010 is shown for visual clarity. (b) The zero-mean residual natural variability component, T_{nat} and the anthropogenic trend, T_{anth} (red and blue, respectively). (c) Spectra of the raw temperature series and the residual component, T_{nat} (blue and red, respectively). The exponent, β was obtained from the linear regression of the smoothed spectrum. The reference dashed line with slope $1+2H$ was also included. (d) Similar to (c), but now considering the average spectra for all the 10512 grid points.

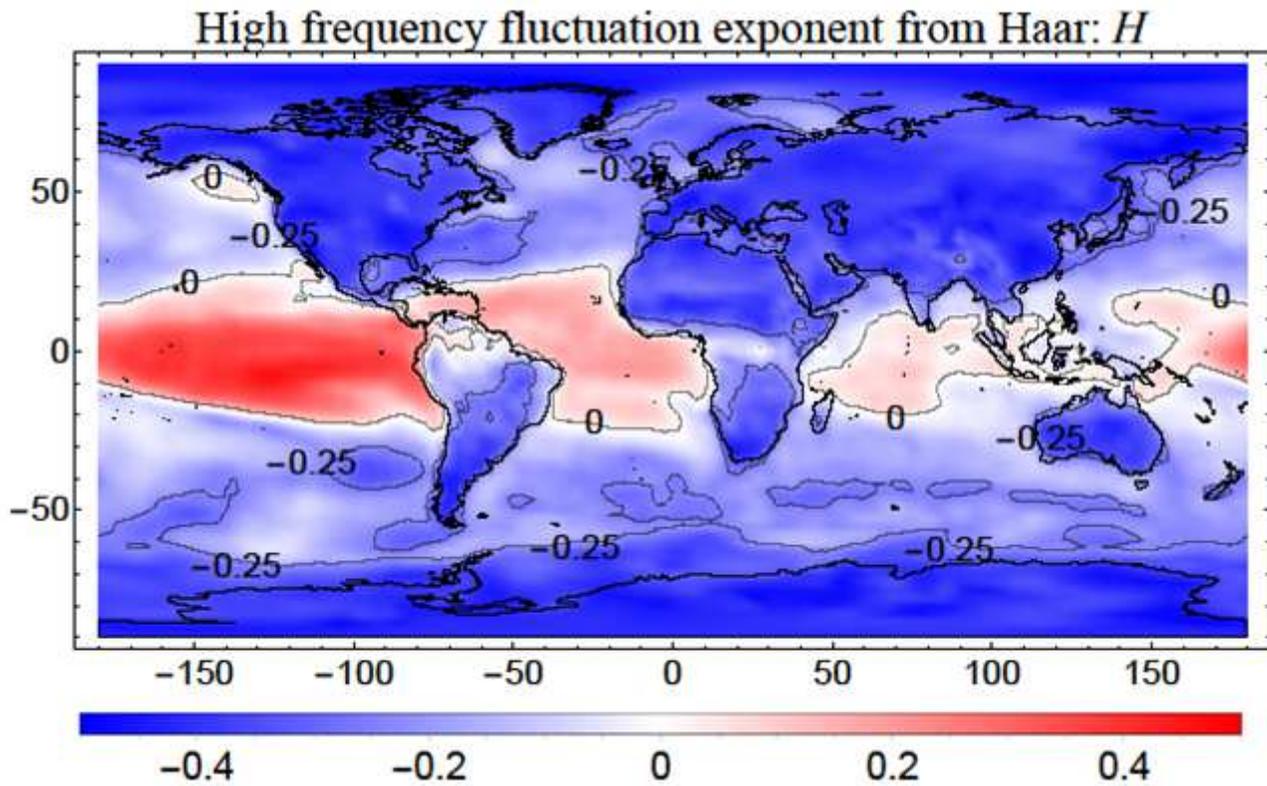


Figure 2

Map of the fluctuation exponents obtained from the Haar fluctuation analysis (Lovejoy and Schertzer 2012a), in the high-frequency scaling regime between 2 months and 2 years. Note: The designations employed and the presentation of the material on this map do not imply the expression of any opinion whatsoever on the part of Research Square concerning the legal status of any country, territory, city or area or of its authorities, or concerning the delimitation of its frontiers or boundaries. This map has been provided by the authors.

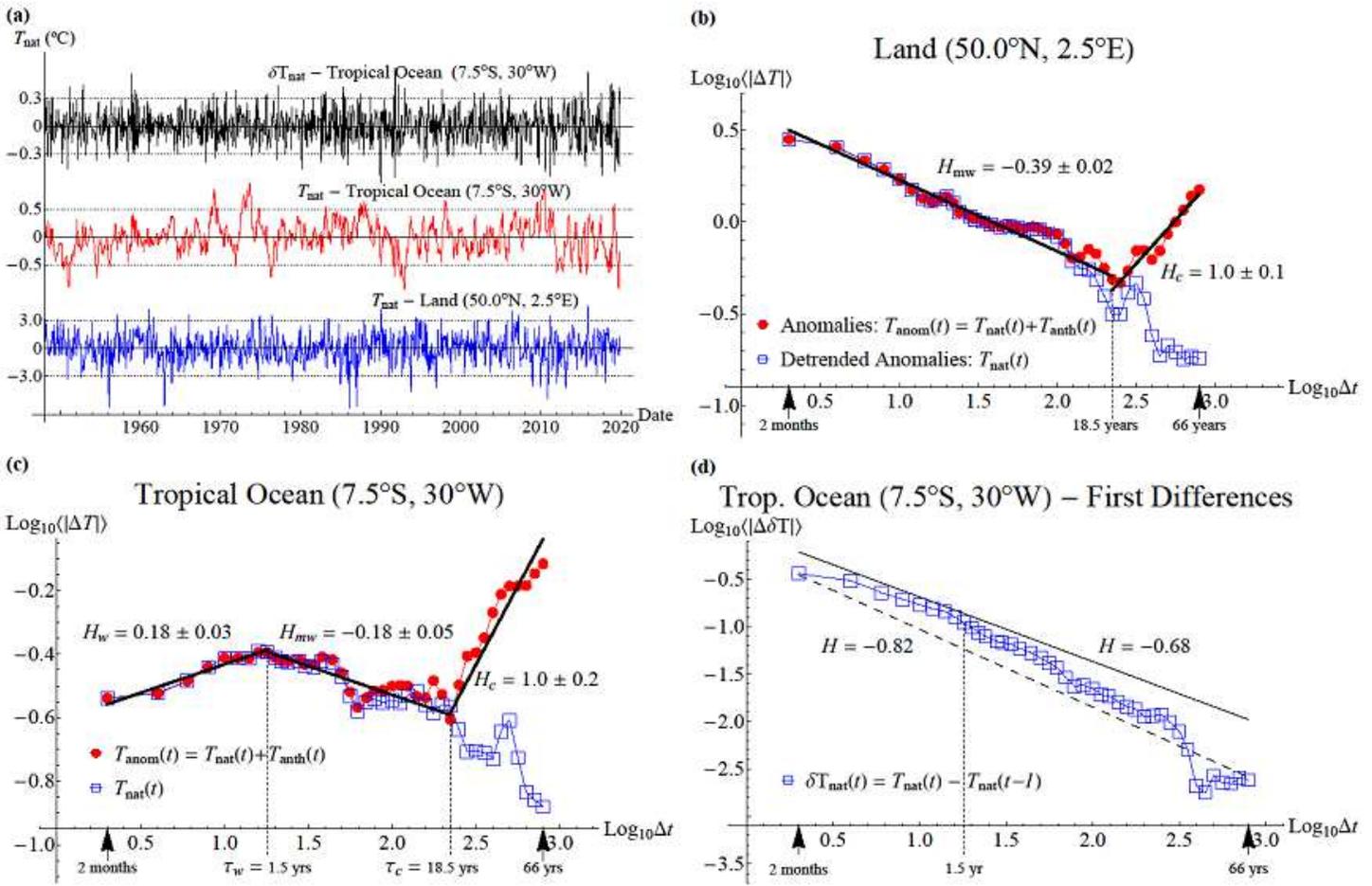


Figure 3

Examples of Haar fluctuation analysis for two points, one over land and one over ocean. (a) In blue, time series for a point over land with coordinates 50.0°N, 2.5°E (same grid point used before in Sect. 2.1); in red, for a point over ocean located at 7.5°S, 30°W and in black, the series of the temperature differences, $\delta T_{nat}(t) = T_{nat}(t) - T_{nat}(t-1)$, for the same point over ocean (increments of the time series in red). (b) Average fluctuation as a function of the time scale before and after removing the anthropogenic trend for the point over land (red line with circles for the anomalies before removing the anthropogenic component and blue line with empty squares for the detrended anomalies). The reference lines with slopes $H_{mw} = -0.39 \pm 0.02$ and $H_c = 1.0 \pm 0.1$ were obtained from regression of the anomalies' fluctuations in the respective macroweather and climate regimes. (c) Same as in (b) but now for the point over ocean. The three regimes (weather, macroweather and climate) are observed for this point. The corresponding transition scales and the respective exponents obtained from linear regression are also included in the graph. (d) Haar fluctuation analysis of the series of increments $\delta T_{nat}(t) = T_{nat}(t) - T_{nat}(t-1)$ for the point over ocean. The dashed line included as reference has slope $H = H_w - 1 = -0.82$, where H_w is the one shown in (c) and the solid line has a slope $H = -0.68$, which is the exponent obtained from the maximum likelihood method assuming that δT_{nat} is a fractional Gaussian noise (fGn) process (see next section).

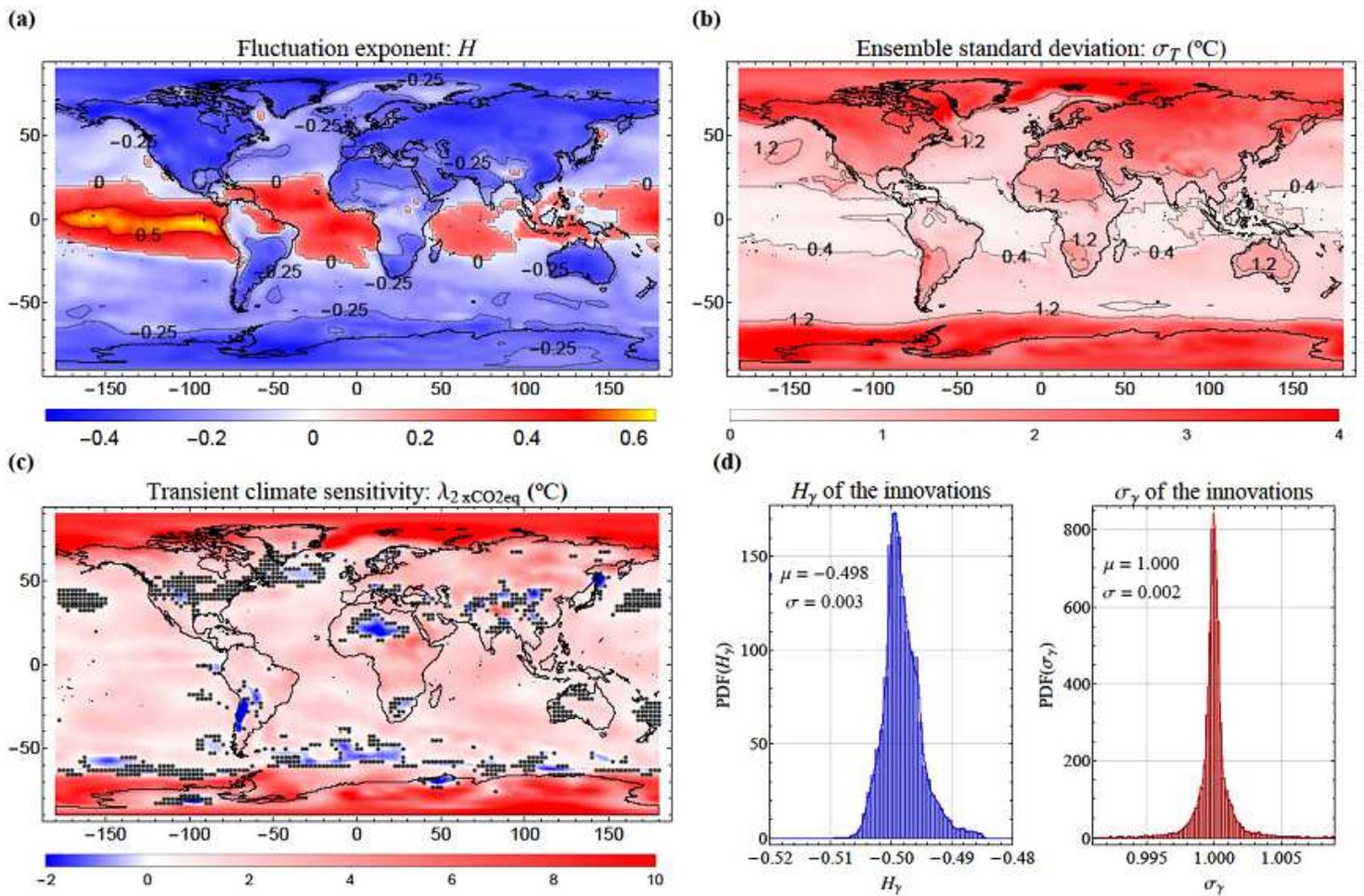


Figure 4

Estimates of the three parameters ($H, \sigma_T, \lambda_{2\times\text{CO}_2\text{eq}}$) obtained for each grid point and statistics of the innovations, $\gamma(t)$. (a) Maximum likelihood estimates of the temperature fluctuation exponent (compare with the estimates shown in Fig. 2). There is a discontinuity from negative to positive values of H as we approach the tropical ocean, corresponding to the change in model from fGn to fBm . (b) The standard deviation, σ_T , of the infinite ensemble fGn process. (c) Map of the transient climate sensitivity, defined in Eq. (2). The places marked with "x" indicate pixels where the null hypothesis, $\lambda_{2\times\text{CO}_2\text{eq}}=0$, cannot be rejected with more than 90% confidence. (d) Histograms of the fluctuation exponent and the standard deviation of the innovations (H_γ and σ_γ , respectively) for the 10512 grid points. From the histograms, we can conclude that the innovations are very close to white noise for the whole planet ($H_\gamma = -0.498 \pm 0.003$ and $\sigma_\gamma = 1.000 \pm 0.002$). Note: The designations employed and the presentation of the material on this map do not imply the expression of any opinion whatsoever on the part of Research Square concerning the legal status of any country, territory, city or area or of its authorities, or concerning the delimitation of its frontiers or boundaries. This map has been provided by the authors.

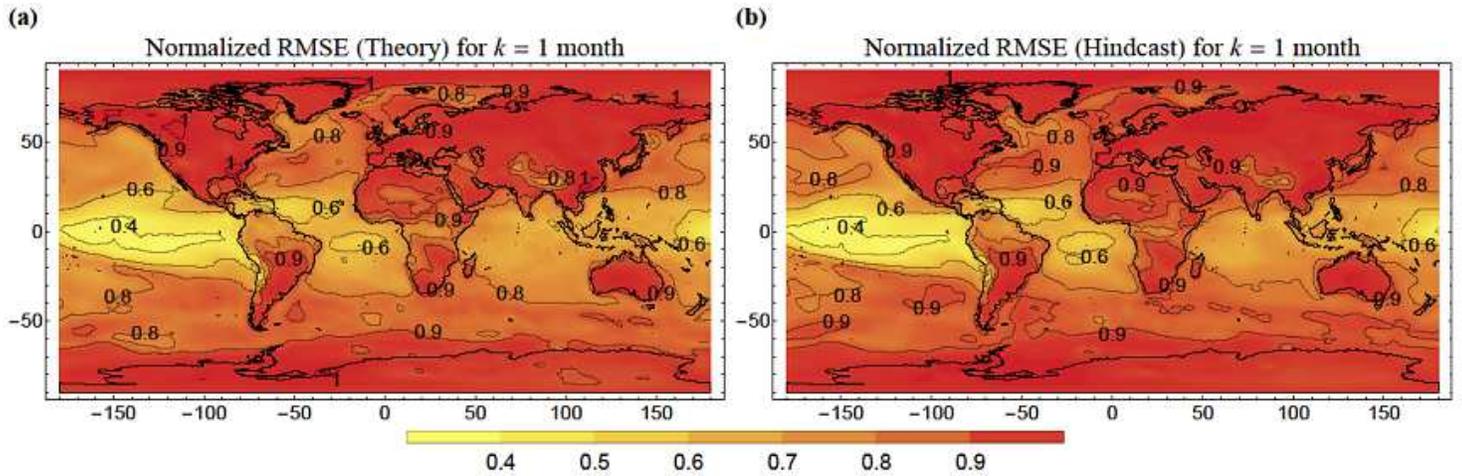


Figure 5

Theoretical and hindcasts NRMSE for $k=1$ month. The corresponding RMSEs were obtained using Eqs. (A13) and (B3), respectively, and the normalization standard deviation from Eq. (5) for the natural variability. Note: The designations employed and the presentation of the material on this map do not imply the expression of any opinion whatsoever on the part of Research Square concerning the legal status of any country, territory, city or area or of its authorities, or concerning the delimitation of its frontiers or boundaries. This map has been provided by the authors.

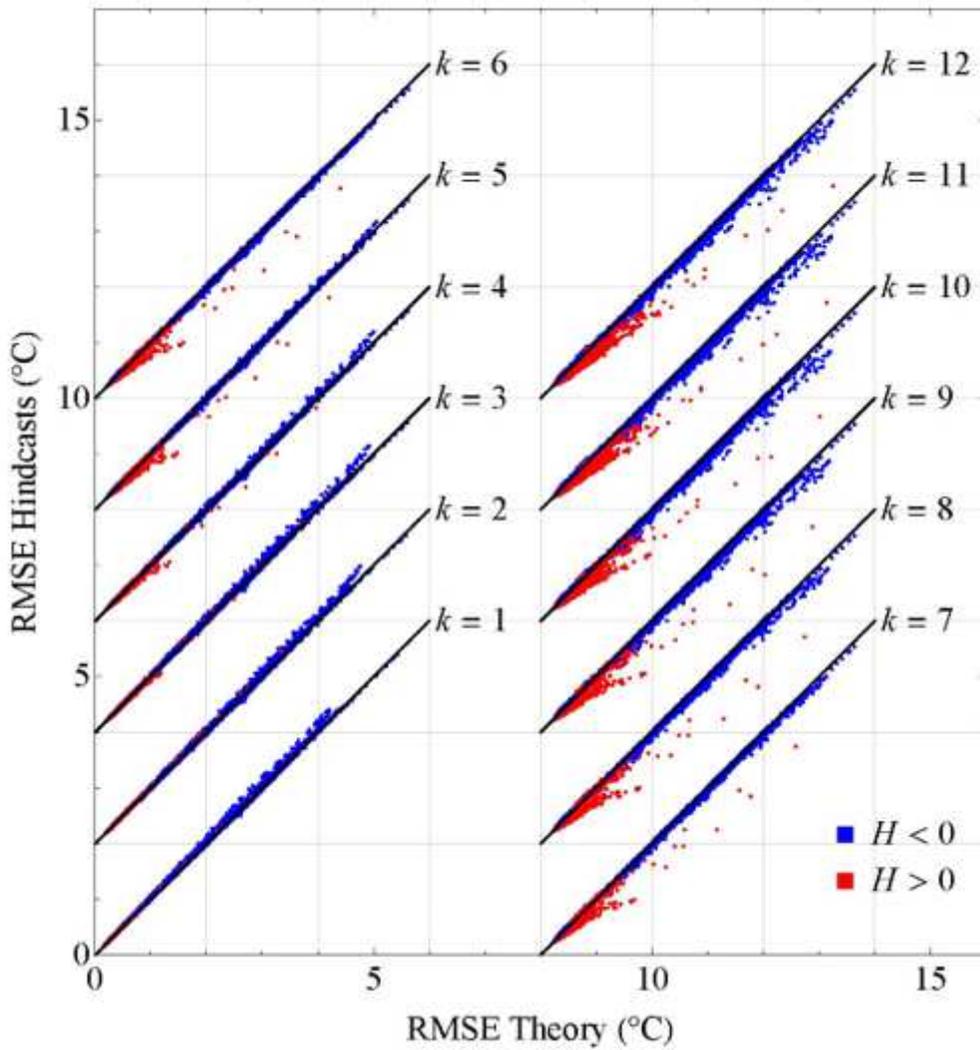


Figure 6

Scatter plots for each horizon including the 10512 grid points, showing the verification RMSE obtained from hindcasts vs. the expected theoretical $RMSE_{nat}^{theory}$ predicted by Eq. (A13). The graphs were displaced vertically by $2^{\circ}C$ (plus a horizontal displacement of $8^{\circ}C$ for $k \geq 7$ months) for visual clarity. The black line at 45° is a reference indicating perfect agreement between theory and verification results. The blue points represent locations where $H < 0$ and the natural variability is modeled as an fGn process and the red points are for places where $H > 0$ and we use the fBm model.

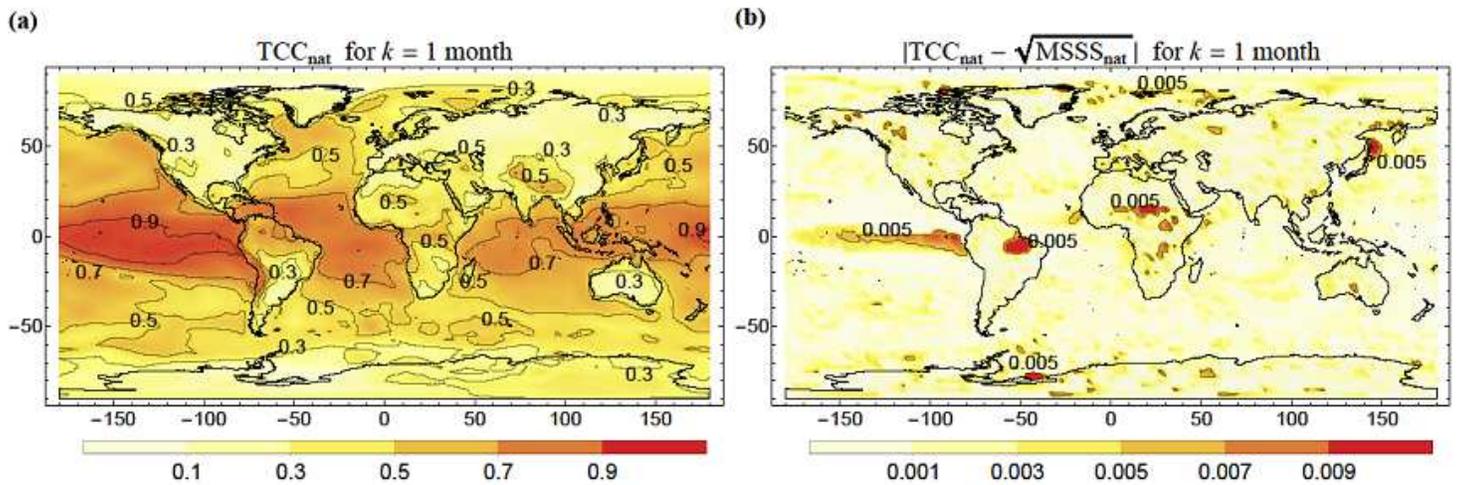


Figure 7

Maps of TCC_{nat} and the absolute difference $|TCC_{nat} - \sqrt{MSSS_{nat}}|$ obtained from hindcasts for $k=1$ month. The colour scale in (b) was rescaled 100 times with respect to (a) so the differences could be perceptible. Note: The designations employed and the presentation of the material on this map do not imply the expression of any opinion whatsoever on the part of Research Square concerning the legal status of any country, territory, city or area or of its authorities, or concerning the delimitation of its frontiers or boundaries. This map has been provided by the authors.

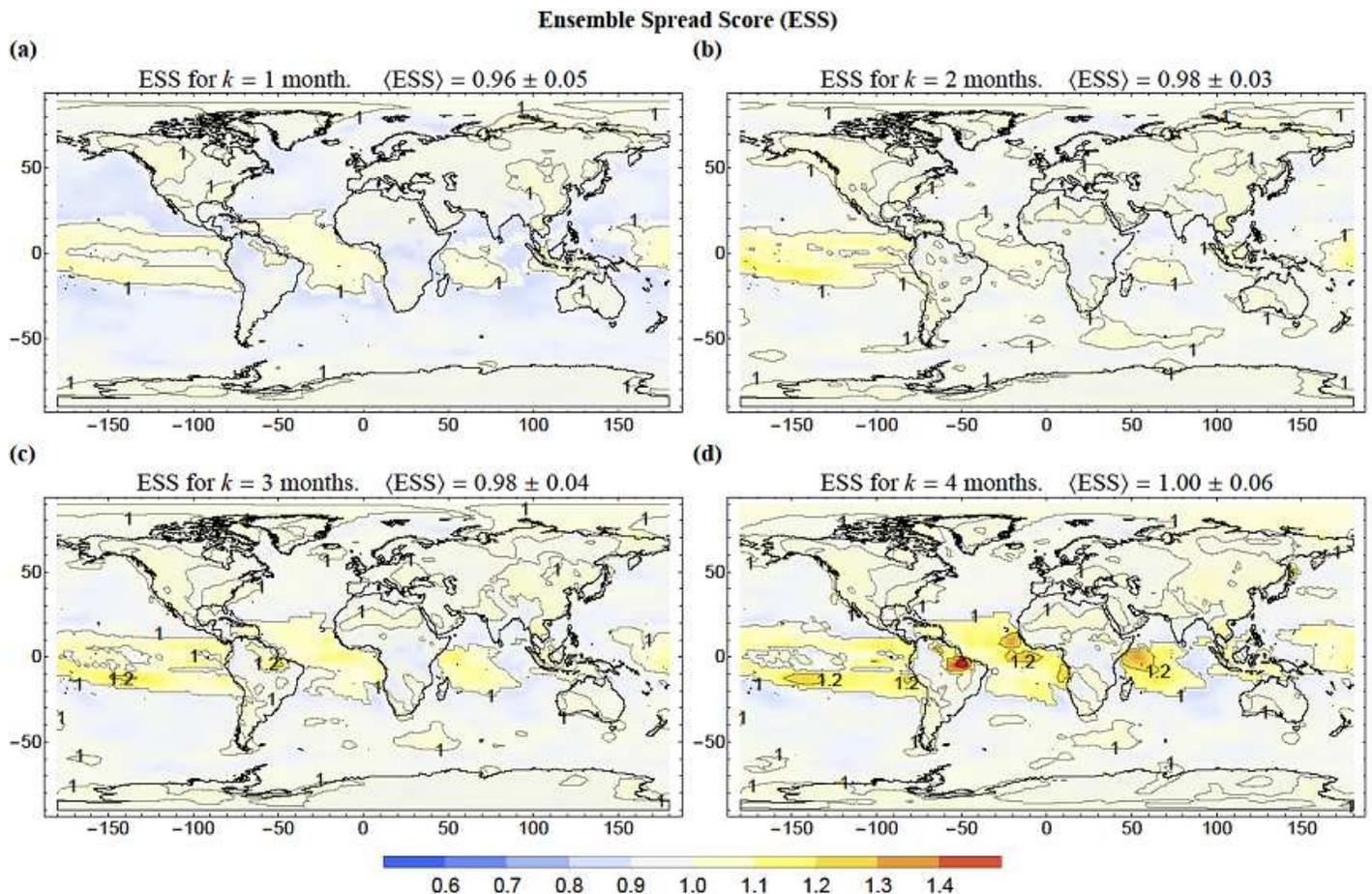


Figure 8

Maps of ESS of StocSIPS for horizons k from 1 to 4 months (panels (a) to (d), respectively). The values of the ESS are very close to 1, with the exception of the tropical ocean where it tends to be “overdispersive” ($ESS > 1$). The average values for the globe with one standard deviation are shown in brackets in the map labels. Note: The designations employed and the presentation of the material on this map do not imply the expression of any opinion whatsoever on the part of Research Square concerning the legal status of any country, territory, city or area or of its authorities, or concerning the delimitation of its frontiers or boundaries. This map has been provided by the authors.

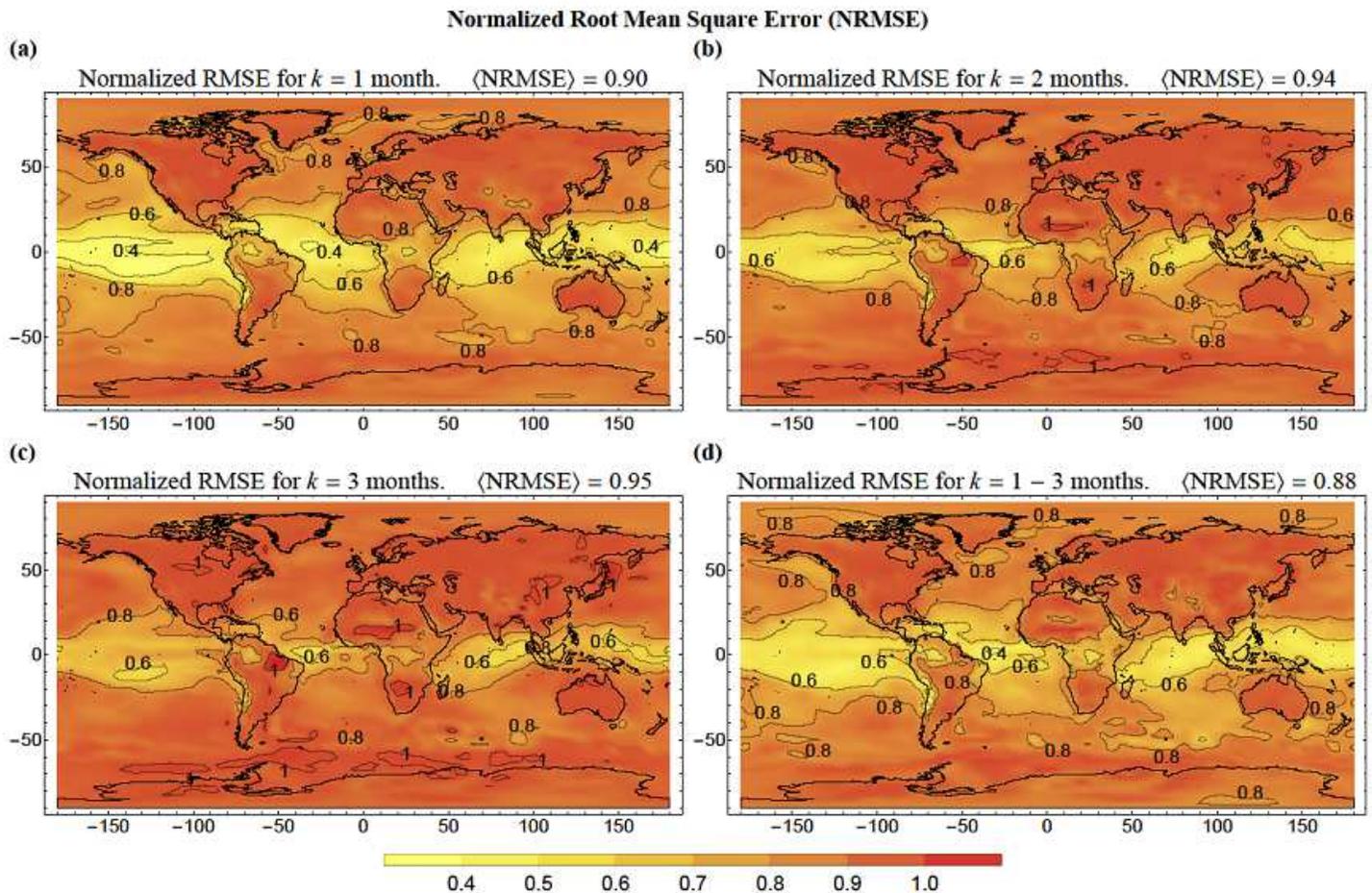


Figure 9

Normalized root mean square error NRMSE for: (a) $k=1$ month, (b) $k=2$ months, (c) $k=3$ months and (d) for the all-seasons mean (average for $k=1-3$ months). The values in brackets in the figure labels represent the areal mean of global NRMSE. Note: The designations employed and the presentation of the material on this map do not imply the expression of any opinion whatsoever on the part of Research Square concerning the legal status of any country, territory, city or area or of its authorities, or concerning the delimitation of its frontiers or boundaries. This map has been provided by the authors.

Mean Square Skill Score (MSSS)

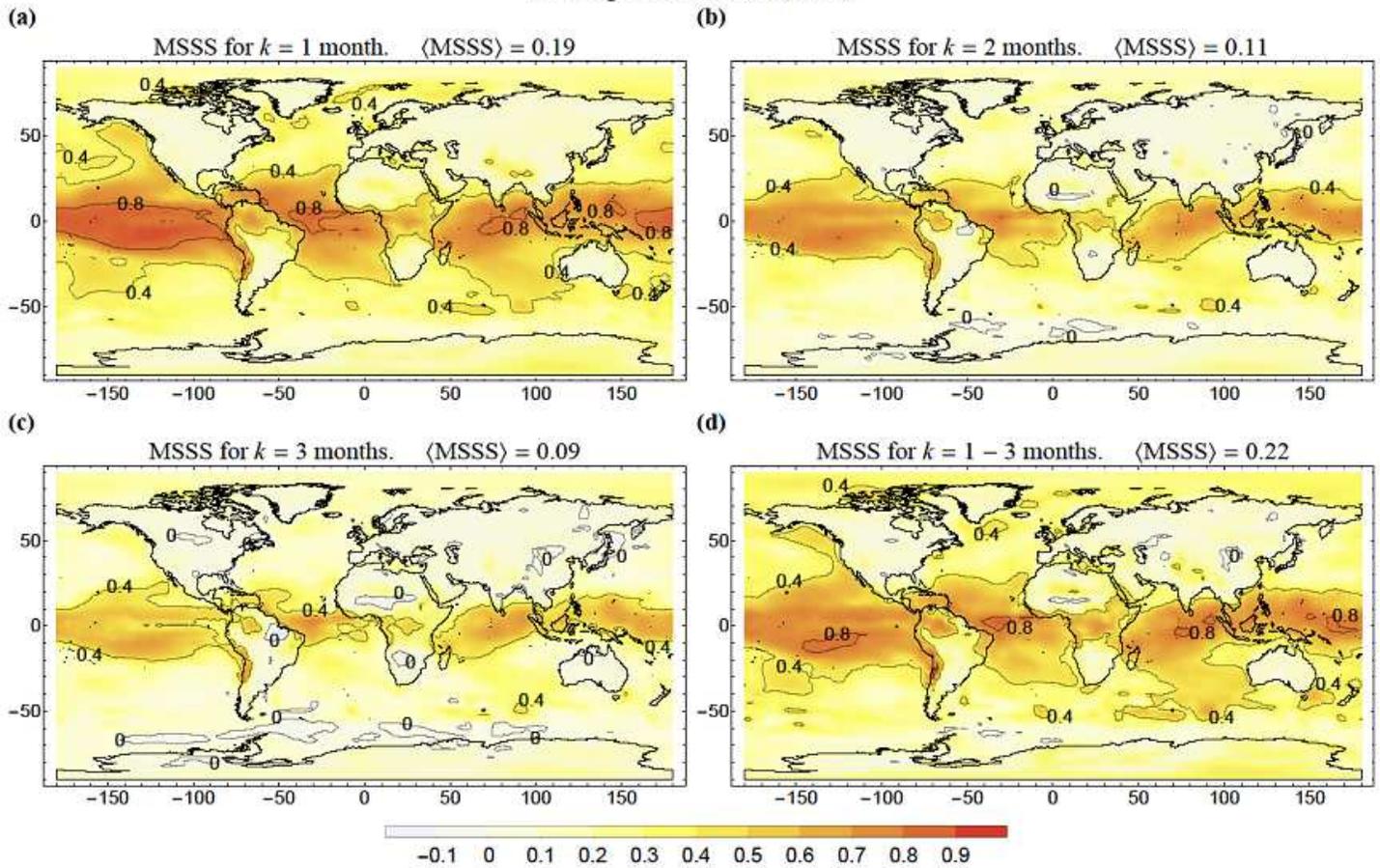


Figure 10

Mean square skill score (MSSS) for: (a) $k=1$ month, (b) $k=2$ months, (c) $k=3$ months and (d) for the all-seasons mean (average for $k=1-3$ months). The values in brackets in the figure labels represent the areal mean of global MSSS. Note: The designations employed and the presentation of the material on this map do not imply the expression of any opinion whatsoever on the part of Research Square concerning the legal status of any country, territory, city or area or of its authorities, or concerning the delimitation of its frontiers or boundaries. This map has been provided by the authors.

Temporal Correlation Coefficient (TCC)

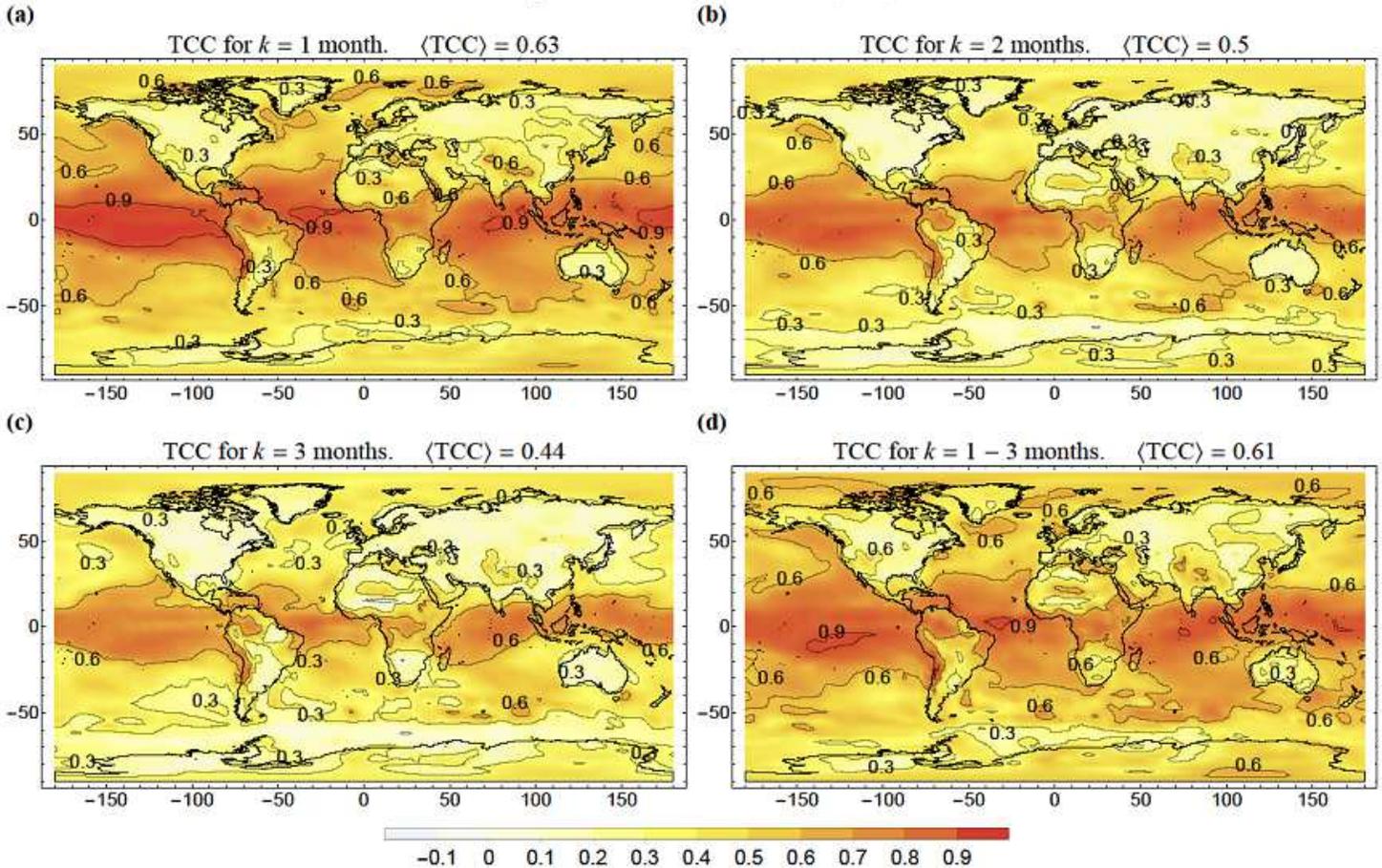


Figure 11

Anomaly correlation coefficient (TCC) for: (a) $k=1$ month, (b) $k=2$ months, (c) $k=3$ months and (d) for the all-seasons mean (average for $k=1-3$ months). The values in brackets in the figure labels represent the areal mean of global TCC. Note: The designations employed and the presentation of the material on this map do not imply the expression of any opinion whatsoever on the part of Research Square concerning the legal status of any country, territory, city or area or of its authorities, or concerning the delimitation of its frontiers or boundaries. This map has been provided by the authors.

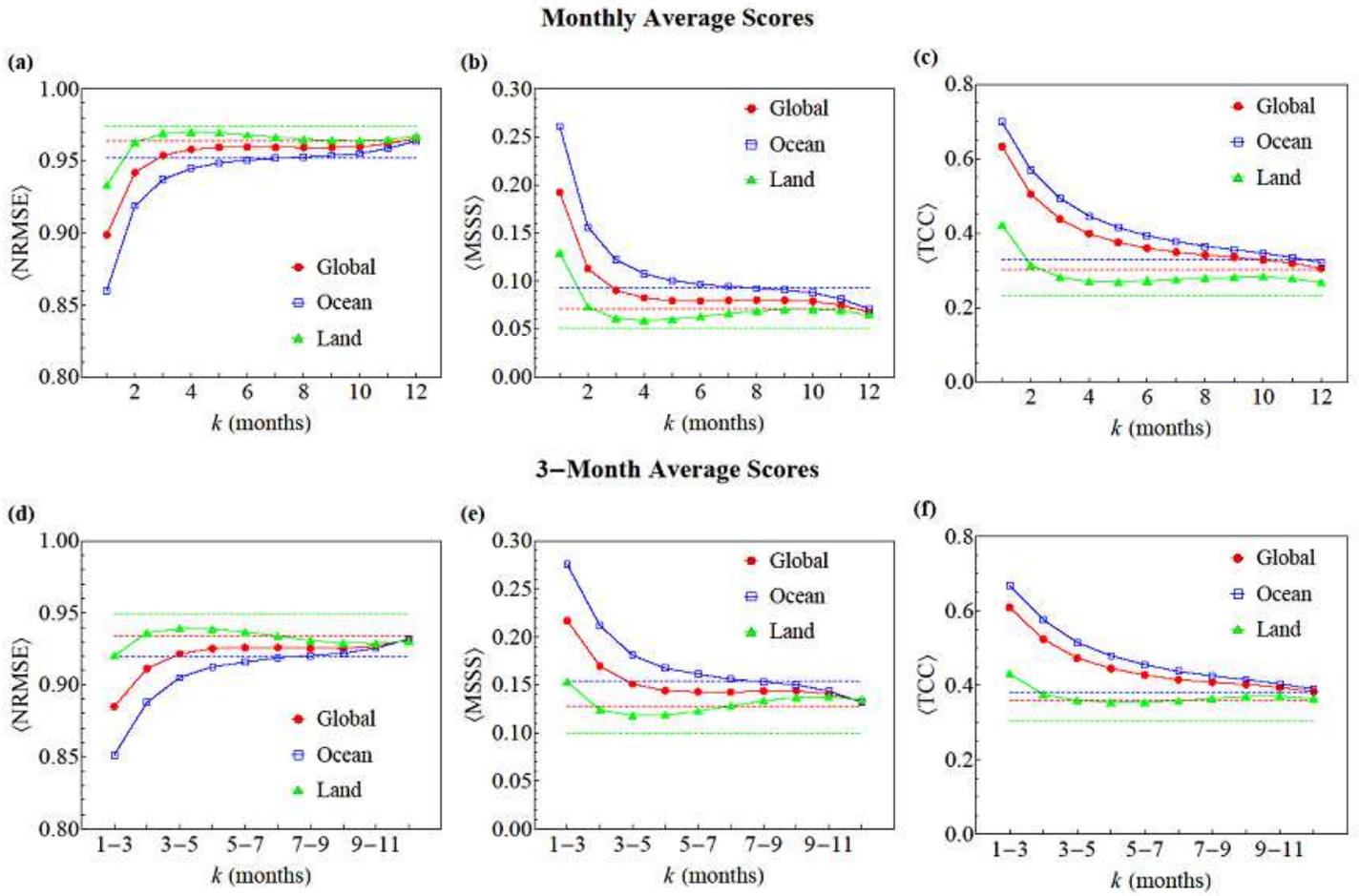


Figure 12

Graphs of the area-averaged NRMSE, MSSS and ACC for the monthly (panels (a), (b) and (c)) and the 3-month average (panels (d), (e) and (f)) forecasts as a function of the forecast horizon. In all the graphs, the red lines with circles correspond to the average considering the grid points for the whole planet, the blue lines with open squares are for places over the ocean and the green lines with triangles are for grid points over land. The corresponding dashed lines of the same colours represent the respective scores obtained if only the anthropogenic trend is forecast.

Seasonal Interannual Variability

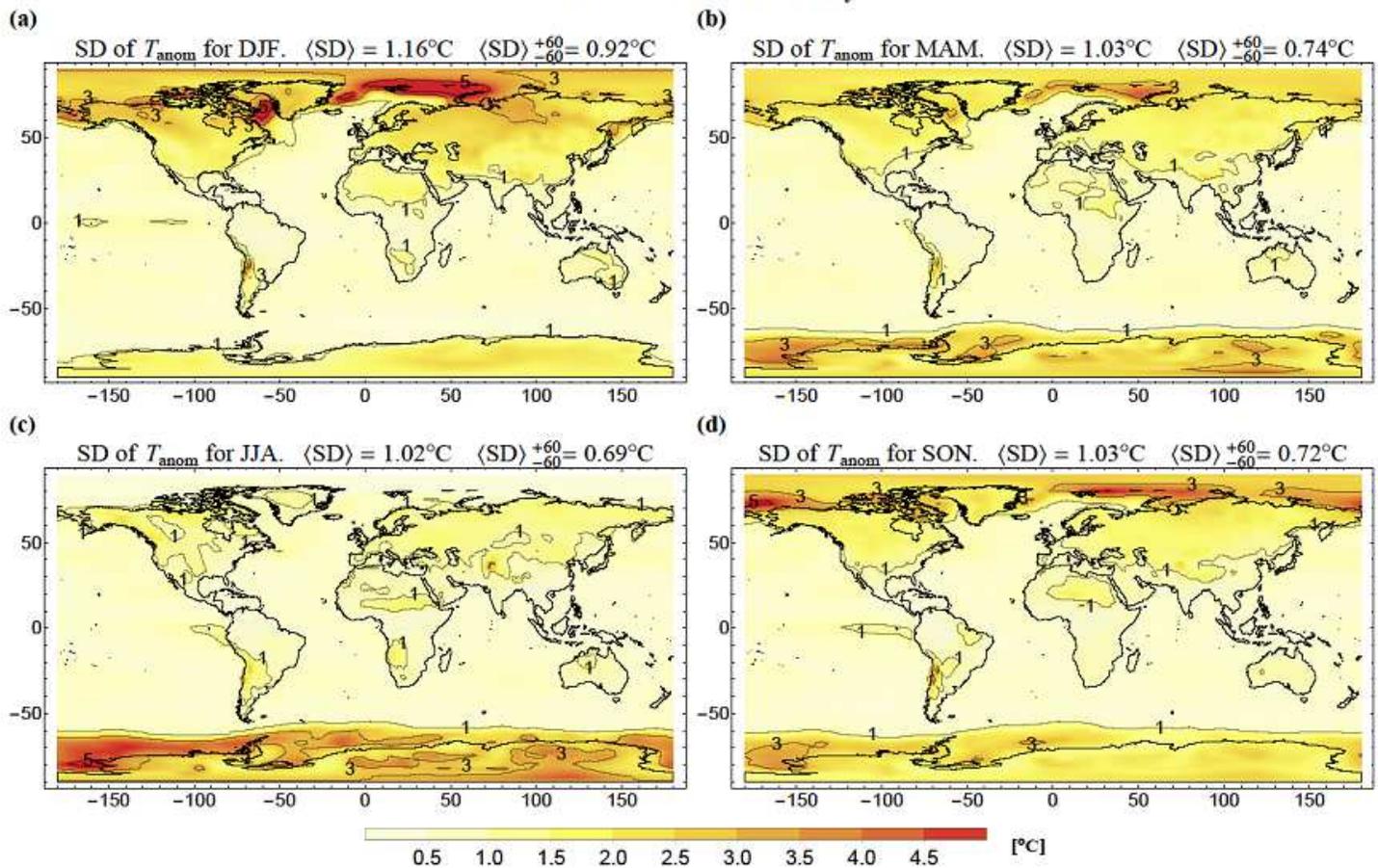


Figure 13

Interannual standard deviation (SD) of the temperature anomalies for the conventional seasons: (a) DJF, (b) MAM, (c) JJA and (d) SON. The values in brackets in the figure labels represent the areal mean of global standard deviation and the areal mean excluding the poles (between 60°S and 60°N). Note: The designations employed and the presentation of the material on this map do not imply the expression of any opinion whatsoever on the part of Research Square concerning the legal status of any country, territory, city or area or of its authorities, or concerning the delimitation of its frontiers or boundaries. This map has been provided by the authors.

Seasonal Mean Square Skill Score (MSSS)

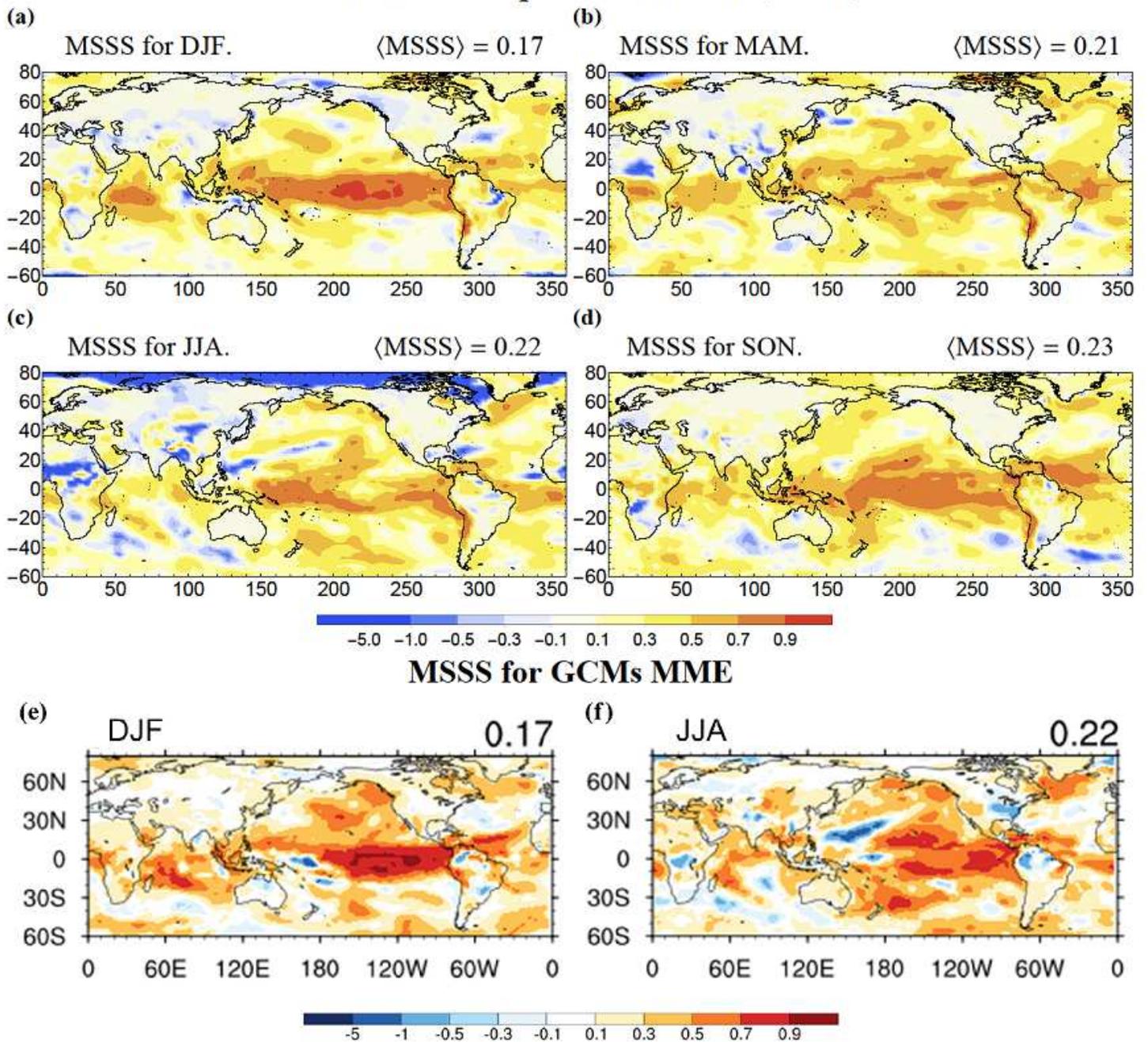


Figure 14

MSSS for: (a) DJF, (b) MAM, (c) JJA and (d) SON. In all cases, the forecasts used data up to the beginning of each respective season (average for $k=1-3$ months). The values in brackets in the figure labels represents the globally averaged MSSS (see Eq. (B10)). The maps shown in panels (e) and (f) for the GCMs MME prediction of DJF and JJA, respectively, were reproduced from Figs. S1 and S2 of (Kim et al. 2020) (supporting information) for their best MME combination method (GA). This is an open access article distributed under the terms of the Creative Commons CC BY license, which permits unrestricted use, distribution, and reproduction in any medium. Note: The designations employed and the presentation of the material on this map do not imply the expression of any opinion whatsoever on the part of

Research Square concerning the legal status of any country, territory, city or area or of its authorities, or concerning the delimitation of its frontiers or boundaries. This map has been provided by the authors.

Seasonal Temporal Correlation Coefficient (TCC)

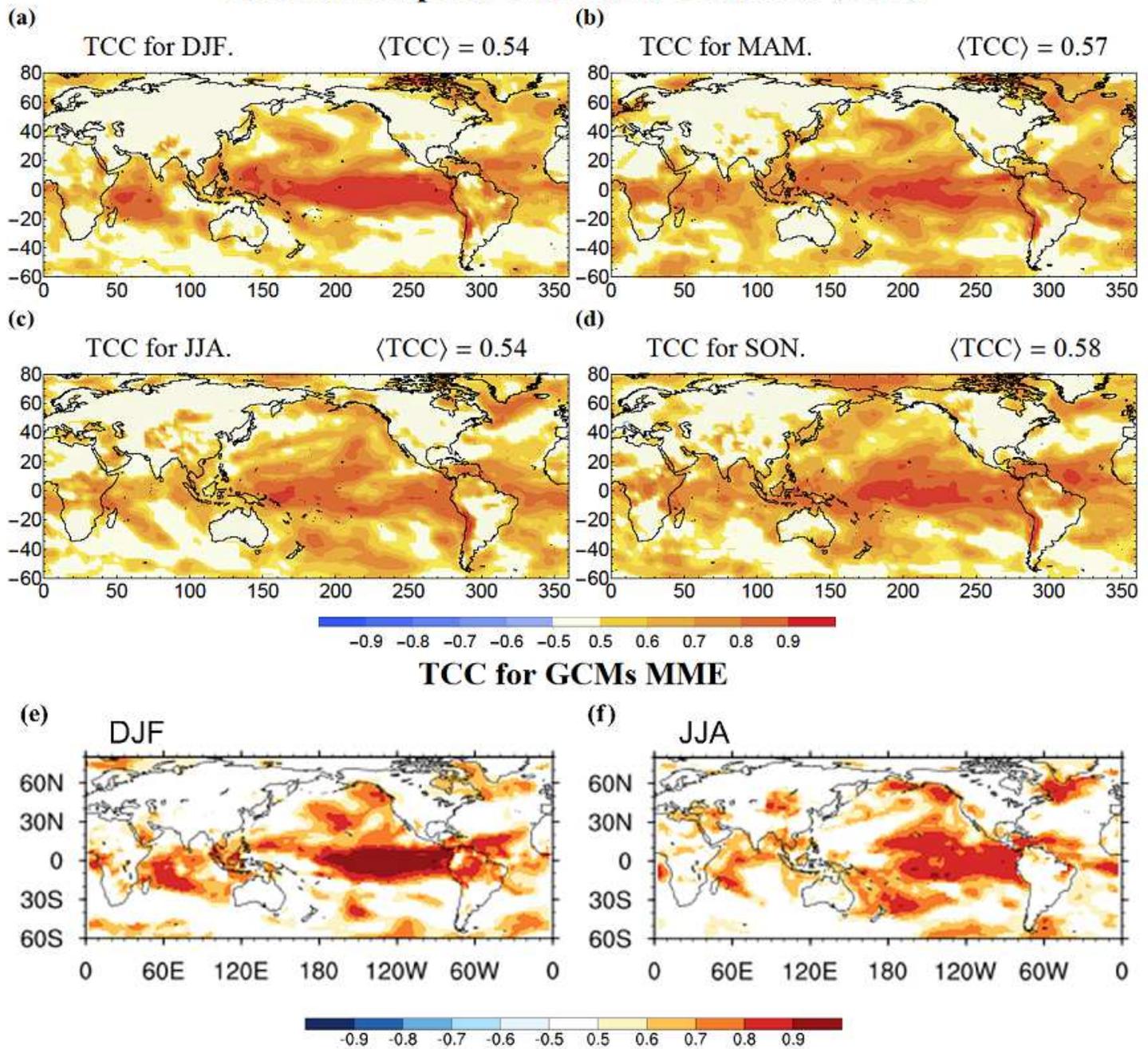


Figure 15

TCC for: (a) DJF, (b) MAM, (c) JJA and (d) SON. In all cases, the forecasts used data up to the beginning of each respective season (average for $k=1-3$ months). The shaded areas indicate the regions over the 5% significance level using two-tailed student's t-test. The values in brackets in the figure labels represent the globally averaged score, $\langle TCC \rangle$, computed using Eq. (B12). The maps shown in panels (e) and (f) for the GCMs MME prediction of DJF and JJA, respectively, were reproduced from Figs. S5 and S6 of (Kim et al. 2020) (supporting information) for their best MME combination method (GA). This is an open access

article distributed under the terms of the Creative Commons CC BY license, which permits unrestricted use, distribution, and reproduction in any medium. Note: The designations employed and the presentation of the material on this map do not imply the expression of any opinion whatsoever on the part of Research Square concerning the legal status of any country, territory, city or area or of its authorities, or concerning the delimitation of its frontiers or boundaries. This map has been provided by the authors.

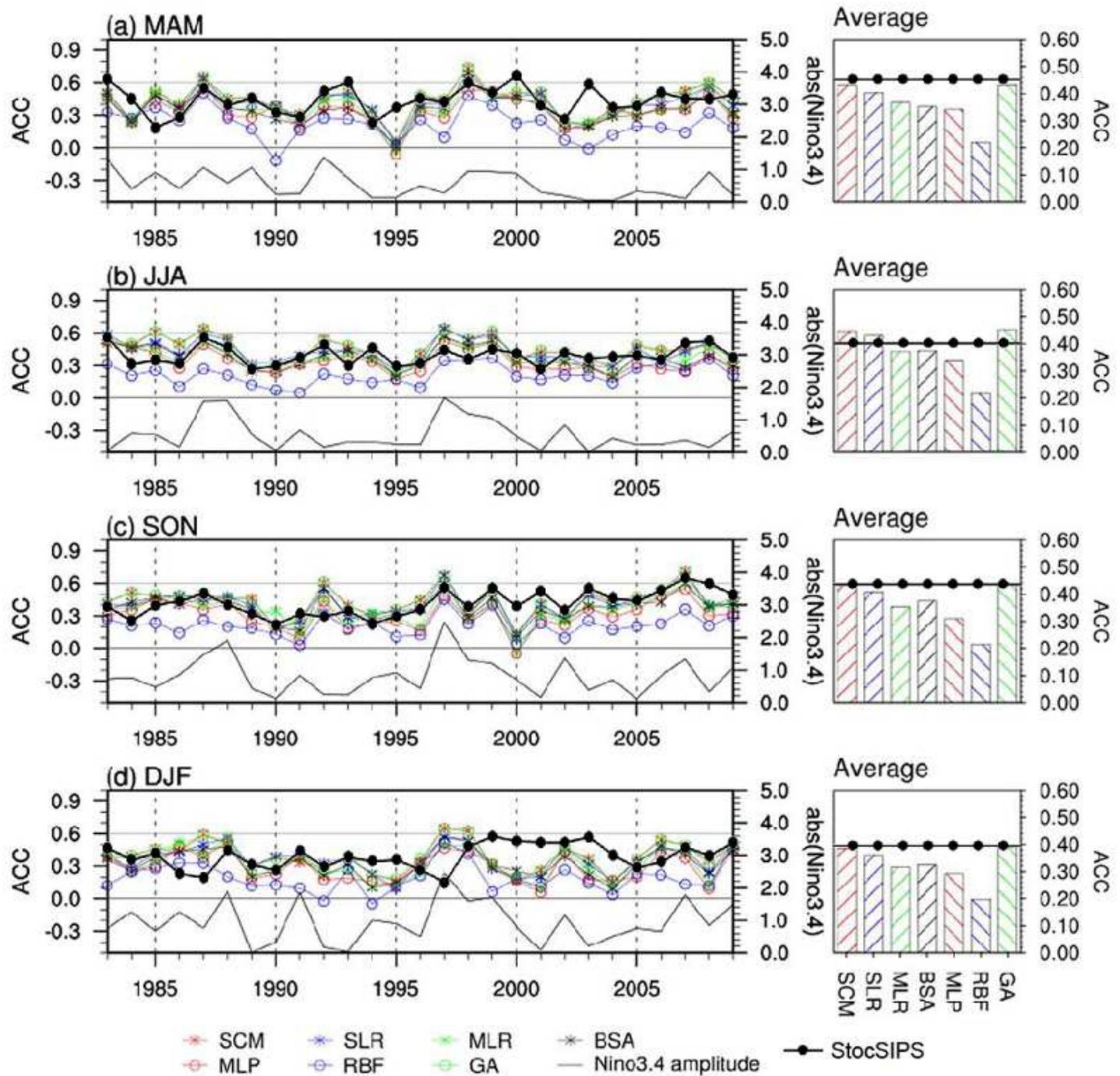


Figure 16

ACC for StocSIPS (black line with solid circles) and for each of the seven MME combination methods studied by Kim et al. (colored lines with markers) in the 27-year verification period 1983-2009 for: (a) MAM, (b) JJA, (c) SON and (d) DJF. The average scores for the POV (see Eq. (B13)) are shown in the right

panels for each of the respective seasons. The absolute value of the El Niño 3.4 index (black line without markers) is also shown. This figure was adapted from Fig. 3 in (Kim et al. 2020) to include the StocSIPS scores. This is an open access article distributed under the terms of the Creative Commons CC BY license, which permits unrestricted use, distribution, and reproduction in any medium.

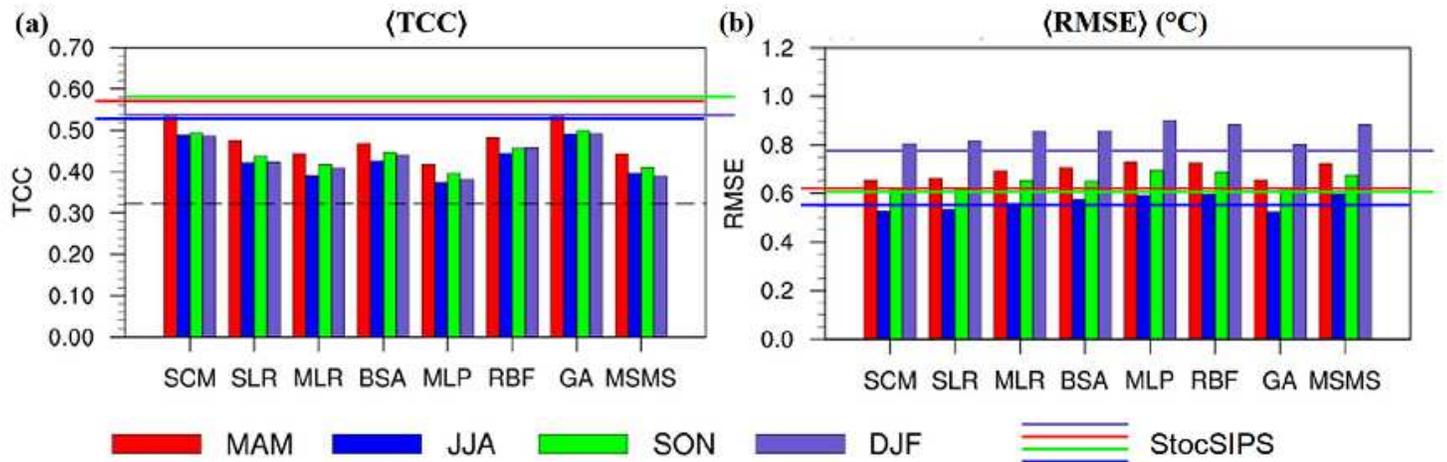


Figure 17

Globally averaged TCC (a) and RMSE (b) (Eqs. (B12) and (B9), respectively) for MAM (red), JJA (blue), SON (green) and DJF (purple) for the period 1983-2009. The bars are for the MME combination methods in (Kim et al. 2020), together with the mean of single model skills (MSMS). The scores for StocSIPS were included as horizontal lines with the same color code for each respective season. The dashed black line indicates that the estimated TCC is statistically significant at the 5% level using the one-tailed Student's t test. This figure was adapted from Figs. 5 and 6 in (Kim et al. 2020) to include the StocSIPS scores. This is an open access article distributed under the terms of the Creative Commons CC BY license, which permits unrestricted use, distribution, and reproduction in any medium.

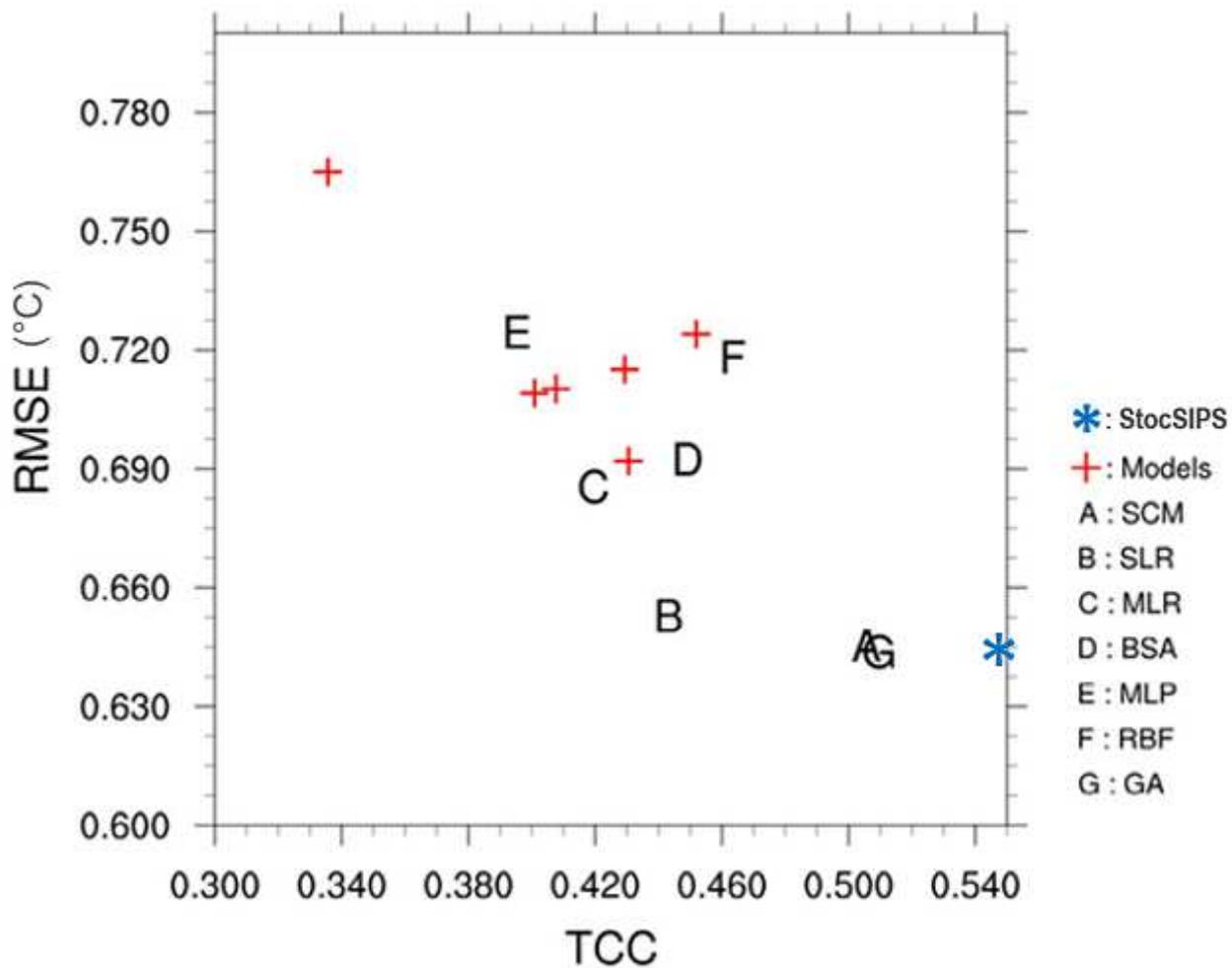


Figure 18

4-season-averaged RMSE vs. TCC for the six individual models used in (Kim et al. 2020) (red crosses), the six MME combinations (letters) and StocSIPS (blue asterisk). This figure was adapted from Figs. 7 in (Kim et al. 2020) to include the StocSIPS scores. This is an open access article distributed under the terms of the Creative Commons CC BY license, which permits unrestricted use, distribution, and reproduction in any medium.

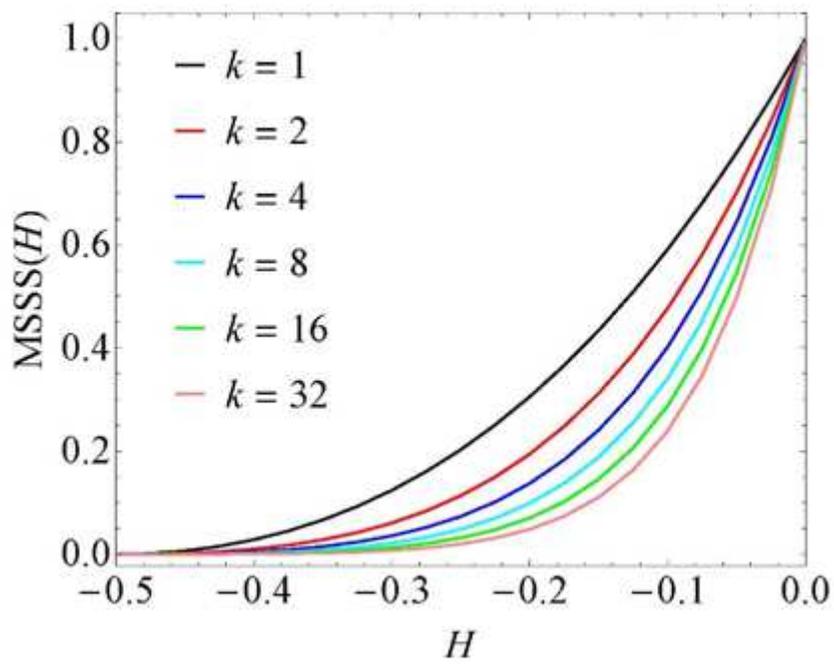


Figure 19

Graphs of the theoretical MSSS (Eq. (A16)) as a function of H for different values of k . A memory $m=50$ was used for computing the MSSS.