

Investigating the Applicability of Emergent Constraints

Alexander J. Winkler^{1,2}, Ranga B. Myneni³, and Victor Brovkin¹

¹Max Planck Institute for Meteorology, Bundesstrasse 53, 20146 Hamburg, Germany

²International Max Planck Research School on Earth System Modelling, Bundesstrasse 53, 20146 Hamburg, Germany

³Department of Earth and Environment, Boston University, Boston, Massachusetts 02215, USA

Correspondence: Alexander J. Winkler (alexander.winkler@mpimet.mpg.de)

1 Abstract.

2 Recent research on Emergent Constraints (EC) has delivered promising results in narrowing down uncertainty in climate pre-
3 dictions. The method utilizes a measurable variable (predictor) from the recent historical past to obtain a constrained estimate
4 of change in an entity of interest (predictand) at a potential future CO₂ concentration (forcing) from multi-model projections.
5 This procedure critically depends on, first, accurate estimation of the predictor from observations and models, and second, on
6 a robust relationship between inter-model variations in the predictor-predictand space. Here, we investigate issues related to
7 these two themes in a carbon cycle case study using observed vegetation greening sensitivity to CO₂ forcing as a predictor
8 of change in photosynthesis (Gross Primary Productivity, GPP) for a doubling of pre-industrial CO₂ concentration. Greening
9 sensitivity is defined as changes in annual maximum of green leaf area index (LAI_{max}) per unit CO₂ forcing realized through
10 its radiative and fertilization effects. We first address the question of how to realistically characterize the predictor of a large
11 area (e.g. greening sensitivity in the northern high latitudes region) from pixel-level data. This requires an investigation into
12 uncertainties in the observational data source and an evaluation of the spatial and temporal variability in the predictor in both
13 the data and model simulations. Second, the predictor-predictand relationship across the model ensemble depends on a strong
14 coupling between the two variables, i.e. simultaneous changes in GPP and LAI_{max}. This coupling depends in a complex man-
15 ner on the magnitude (level), time-rate of application (scenarios) and effects (radiative and/or fertilization) of CO₂ forcing. We
16 investigate how each one of these three aspects of forcing can impair the EC estimate of the predictand (Δ GPP). Our results
17 show that uncertainties in the EC method ~~can~~ primarily originate from a lack of predictor comparability between models and
18 observations, temporal variability, and the observational data source of the predictor. The disagreement between models on
19 the mechanistic behavior of the system under intensifying forcing limits the EC applicability. The here illustrated limitations
20 and sources of uncertainty in the EC method go beyond carbon cycle research and are generally applicable in Earth system
21 sciences.

22 Copyright statement.

of what predictor and predictand?
affect
reword

1 1 Introduction

2 Earth system models (ESMs) are powerful tools to predict responses to a variety of forcings such as increasing atmospheric
3 concentration of greenhouse gases and other agents of radiative forcing (Klein and Hall, 2015). Still, longterm ESM projections
4 of climate change have substantial uncertainties. This can be due to poorly understood processes in some cases, and in others,
5 to missing or simplified representations called parameterizations (Flato et al., 2013; Klein and Hall, 2015; Knutti et al., 2017).
6 Certain important processes, especially in the atmosphere, happen at spatial scales finer than can be possibly represented in
7 current ESMs. Consequently, various phenomena in the system ranging from local extreme precipitation events to large-scale
8 climate modes, can be poorly simulated (Flato et al., 2013). Errors propagate and can be amplified through feedbacks among
9 interacting components in the Earth system, resulting in biases whose origins can be difficult to identify (Flato et al., 2013).
10 Furthermore, an inherent component of the Earth climatic system, its internal natural variability, is complicated to represent
11 and simulate in models (Flato et al., 2013; Klein and Hall, 2015).

12
13 Model Intercomparison Projects explore these uncertainties by coordinating a wide range of simulation setups focusing on
14 internal variability, boundary conditions, parameterizations, etc. (Taylor et al., 2012; Flato et al., 2013; Eyring et al., 2016;
15 Knutti et al., 2017). Models developed at various institutions are driven with the same forcing information (e.g. historical forc-
16 ing) or with identical idealized boundary conditions. However, each modeling group decides which of the processes to consider
17 and implement in their ESM. The conventional approach of handling these multi-model ensembles is to use unweighted ensem-
18 ble averages (Knutti, 2010; Knutti et al., 2017). This assumes that the models are independent of one another and equally good
19 at simulating the climate system (Flato et al., 2013; Knutti et al., 2017). The large spread between model projections suggests
20 that this assumption is not valid. Therefore, alternate methods have been developed to extract results more accurate than multi-
21 model averages (e.g. model weighting scheme based on performance and interdependence, Knutti et al., 2017). The concept of
22 *Emergent Constraints* arises in this context, namely, as a method to reduce uncertainty in ESM projections relying on histori-
23 cal simulations and observations (Hall and Qu, 2006; Boé et al., 2009; Cox et al., 2013; Klein and Hall, 2015; Cox et al., 2018).

24
25 The two key parts of an Emergent Constraint (EC) based method are a linear relationship arising from the collective behavior
26 of a multi-model ensemble and an observational estimate for imposing the said constraint (Fig. 1). The linear relationship is a
27 physically (or physiologically) based correlation between inter-model variations in an observable entity of the contemporary
28 climate system (*predictor*) and a projected variable (*predictand*) that is difficult to observe or not observable at all. Combining
29 the emergent linear relationship with observations of the predictor sets a constraint on the predictand (Cox et al., 2013; Flato
30 et al., 2013; Klein and Hall, 2015; Knutti et al., 2017). Many such ECs have been identified and reported, as briefly summarized
31 below.

32
33 Hall and Qu (2006) proposed a constraint on projections of snow-albedo feedback based on the correlation between large
34 inter-model variations in feedback strength of the current seasonal cycle. The EC was first established for the CMIP3 ensemble

1 and confirmed for phase five of the Coupled Model Intercomparison Project (CMIP5; Flato et al., 2013; Qu and Hall, 2014).
2 Several EC studies followed with the goal of reducing uncertainty in projections of the cloud feedback under global warming,
3 as reviewed by Klein and Hall (2015). It is thought that erroneous representation of low-cloud feedback in ESMs contributes
4 essentially to the large uncertainty in equilibrium climate sensitivity (ECS, 1.5 to 5 K), i.e. warming for a doubling of pre-
5 industrial atmospheric CO₂ concentration (2×CO₂; Sherwood et al., 2014; Klein and Hall, 2015). Recently, Cox et al. (2018)
6 presented a different approach to constrain ECS based on its relationship to variability of global temperatures during the recent
7 historical warming period. They reported a constrained ECS estimate of 2.8 K for 2×CO₂ (66% confidence limits of 2.2 – 3.4
8 K).

9
10 The concept of EC also found its way into the field of carbon cycle projections. A series of studies analyzed the extent
11 to which inter-annual atmospheric CO₂ variability can serve as a predictor of longterm temperature sensitivity of terrestrial
12 tropical carbon storage. Cox et al. (2013) and Wenzel et al. (2014) reported an emergent linear relationship, although with
13 different slopes for CMIP3 and CMIP5 ensembles, resulting in slightly divergent constrained estimates (CMIP3: -53 ± 17 Pg
14 C K⁻¹, CMIP5: -44 ± 14 Pg C K⁻¹). Wang et al. (2014) however were unable to detect a similar relationship between the
15 proposed predictor and predictand. Recently, Lian et al. (2018) presented an EC estimate of the global ratio of transpiration
16 to total terrestrial evapotranspiration (T/ET), which is substantially higher (0.62 ± 0.06) than the unconstrained value ($0.41 \pm$
17 0.11). For the marine tropical carbon cycle, Kwiatkowski et al. (2017) identified an emergent relationship between the longterm
18 sensitivity of tropical ocean net primary production (NPP) to rising sea surface temperature (SST) in the equatorial zone and
19 the interannual sensitivity of NPP to El Niño/Southern Oscillation driven SST anomalies. Tropical NPP is projected to decrease
20 by $3 \pm 1\%$ for 1 K increase in equatorial SST according to the observational constraint.

21
22 Similar results were reported for modeled extra-tropical terrestrial carbon fixation in a 2×CO₂ world. Plant productivity is
23 expected to increase due to the fertilizing and radiative effects of rising atmospheric CO₂ concentration. Wenzel et al. (2016)
24 focused on constraining the CO₂ fertilization effect on plant productivity in the northern high latitudes (60° N – 90° N, NHL)
25 and the entire extra-tropical area in the northern hemisphere (30° N – 90° N) using the seasonal amplitude of longterm CO₂
26 measurements at different latitudes. They presented a linear relationship between the sensitivity of CO₂ amplitude to rising
27 atmospheric CO₂ concentration and the relative increase in zonally averaged gross primary production (GPP) for 2×CO₂. The
28 observed CO₂ amplitude sensitivities at respective stations provided a constraint on the strength of the CO₂ fertilization effect,
29 namely $37\% \pm 9\%$ and $32\% \pm 9\%$ for the NHL and the extra-tropical region, respectively.

30 *↳ not sure what these numbers represent*

31 Focusing on the NHL, Winkler et al. (2019) investigated how both effects of CO₂ enhance plant productivity while assess-
32 ing the feasibility of vegetation greenness changes as a constraint. Enhanced GPP due to the physiological effect and ensuing
33 climate warming is indirectly evident in large-scale increase in summer time green leaf area (Myneni et al., 1997a; Zhu et al.,
34 2016). Historical CMIP5 simulations show that the maximum annual leaf area index (LAI_{max}, leaf area per ground area) in-
35 creases linearly with both CO₂ concentration and temperature in NHL. In all ESMs, these changes in LAI_{max} strongly correlate

1 to changes in GPP arising from the combined radiative and physiological effects of CO₂ enrichment. Thus, the large variation
2 in modeled historical LAI_{max} responses to the effects of CO₂ linearly maps to variation in ΔGPP at 2×CO₂ in the CMIP5
3 ensemble. This linear relationship in inter-model variations enables the usage of the observed longterm change in LAI_{max} as
4 an EC on ΔGPP at 2×CO₂ in NHL ($3.4 \pm 0.2 \text{ Pg C yr}^{-1}$; Winkler et al., 2019).

Is this the increase in GPP at 2×CO₂ per unit-increase in LAI?

5
6 The robustness of these EC estimates is debated, mainly because the EC approach is susceptible to methodological incon-
7 sistencies. For example, Cox et al. (2013), Wang et al. (2014) and Wenzel et al. (2015) investigated on constraining future
8 terrestrial tropical carbon storage using the same set of models and data. However, they arrived at different EC estimates and
9 divergent conclusions. Some reasons for failure and essential criteria of the EC approach were described previously (Bracegir-
10 dle and Stephenson, 2012b; Klein and Hall, 2015), but this list is far from complete. To account for this gap in the literature,
11 a detailed investigation and description of the EC method in terms of its potential sources of uncertainty and the range of
12 applicability are needed.

13
14 Here, we revisit the study of Winkler et al. (2019) and elaborate on key issues concerning the robustness of the EC method.
15 Uncertainty of the constrained estimate depends on (a) observed predictor and (b) modeled relationship, aside from the
16 goodness-of-fit of the latter (green shading in Fig. 1). As for (a), the source of observations is an obvious first line of in-
17 quiry (Sect. 3.1). Spatial aggregation of data and model simulations introduces uncertainties, as the EC method is applied on
18 large areal values of predictor and predictand. This is the subject of Sect. 3.2. The observed and modeled predictors are from
19 the historical period. The representativeness, duration and match between data and models all introduce an uncertainty related
20 to variations in the temporal domain – these are explored in Sect. 3.3. The yellow shading in Fig. 1 represents the total uncer-
21 tainty on observed predictor from these three fronts. Regarding (b), the modeled linear relation varies (grey shading in Fig. 1)
22 depending on three attributes of the forcing, i.e. CO₂ concentration change, its magnitude, rate and effect (Sect. 3.4 and 3.5).
23 Lessons learned from analyses along these lines are presented in the conclusion section at the end.

24

1 2 Data and Methods

2 2.1 Remotely sensed leaf area index

3 We used the recently updated version (V1) of the leaf area index dataset (LAI3g) developed by (Zhu et al., 2013). It was gen-
4 erated using an artificial neural network (ANN) and the latest version (third generation) of the Global Inventory Modeling and
5 Mapping Studies group (GIMMS) Advanced Very High Resolution Radiometer (AVHRR) normalized difference vegetation
6 index (NDVI) data (NDVI3g). The latter have been corrected for sensor degradation, inter-sensor differences, cloud cover, ob-
7 servational geometry effects due to satellite drift, Rayleigh scattering and stratospheric volcanic aerosols (Pinzon and Tucker,
8 2014). This dataset provides global and year-round LAI observations at 15-day (bi-monthly) temporal resolution and 1/12
9 degree spatial resolution from July 1981 to December 2016. Currently, this is the only available record of such length.

10

11 The quality of previous version (V0) of LAI3g dataset was evaluated through direct comparisons with ground measurements
12 of LAI and indirectly with other satellite-data based LAI products, and also through statistical analysis with climatic variables,
13 such as temperature and precipitation variability (Zhu et al., 2013). The LAI3gV0 dataset (and related fraction vegetation-
14 absorbed photosynthetically active radiation dataset) has been widely used in various studies (Anav et al., 2013; Piao et al.,
15 2014; Poulter et al., 2014; Forkel et al., 2016; Zhu et al., 2016; Mao et al., 2016; Mahowald et al., 2016; Keenan et al., 2016).
16 The new version, LAI3gV1, used in our study is an update of that earlier version.

17

18 We also utilized a more reliable but shorter dataset from the Moderate Resolution Imaging Spectroradiometer (MODIS)
19 aboard the NASA's Terra satellite (Yan et al., 2016a, b). These data are well calibrated, cloud-screened and corrected for at-
20 mospheric effects, especially tropospheric aerosols. The sensor-platform is regularly adjusted to maintain a precise orbit. All
21 algorithms, including the LAI algorithm, are physics-based, well-tested and currently producing sixth generation datasets. The
22 dataset provides global and year-round LAI observations at 16-day (bi-monthly) temporal resolution and 0.05° spatial resolu-
23 tion from 2000 to 2016.

24

25 Leaf area index is defined as the one-sided green leaf area per unit ground area in broadleaf canopies and as one-half the
26 green needle surface area in needleleaf canopies in both observational and CMIP5 simulation datasets. It is expressed in units
27 of m² green leaf area per m² ground area. Leaf area changes can be represented either by changes in annual maximum LAI
28 (LAI_{max}; Cook and Pau, 2013), or growing season average LAI. In this study, we use the former because of its ease and
29 unambiguity, as the latter requires quantifying the start- and end-dates of the growing season, something that is difficult to do
30 accurately in NHL (Park et al., 2016) with the low resolution model data. Further, LAI_{max}, is less influenced by cloudiness and
31 noise; accordingly, it is most useful in investigations of long-term greening and browning trends. The drawback of LAI_{max}, is
32 the saturation effect at high LAI values (Myneni et al., 2002). However, this is less of a problem in high latitudinal ecosystems
33 which are less-densely vegetated, with LAI_{max}, values typically in the range of 2 to 3.

34

↙ compared to tropics

Keep notation consistent

1/20 degree

1 The bi-monthly satellite datasets were merged to a monthly temporal resolution by averaging the two composites in the same
2 month and bi-linearly remapped to the resolution of the applied reanalysis product ($0.5^\circ \times 0.5^\circ$, CRU TS4.01).

3

4 **2.2 Environmental driver variables**

5 We use time series of temperature and CO_2 to derive the observed historical forcing (Sect. 2.4) and climatologies of pre-
6 cipitation and temperature to calculate climatic regimes (Fig. 2). Monthly averages of near-surface air temperature and pre-
7 cipitation are from the latest version of the Climatic Research Unit Timeseries dataset (CRU TS4.01). The global data are
8 gridded to $0.5^\circ \times 0.5^\circ$ resolution (Harris et al., 2014). Global monthly means of atmospheric CO_2 concentration are from
9 the GLOBALVIEW-CO2 product (obspack_co2_1_GLOBALVIEWplus_v2.1_2016_09_02; for details see [https://doi.org/10.](https://doi.org/10.25925/20190520)
10 [25925/20190520](https://doi.org/10.25925/20190520)) provided by the National Oceanic and Atmospheric Administration / Earth System Research Laboratory
11 (NOAA / ESRL).

12

13 **2.3 Earth system model simulations**

14 We analyzed recent climate-carbon simulations of seven ESMs participating in the fifth phase of the Coupled Model Inter-
15 comparison Project, CMIP (Taylor et al., 2012). The model simulated data were obtained from the Earth System Grid Federa-
16 tion, ESGF (<https://esgf-data.dkrz.de/projects/esgf-dkrz/>). Seven ESMs provide output for the variables of interest (GPP, CO_2 ,
17 LAI, and near-surface air temperature) for simulations titled esmHistorical, RCP4.5, RCP8.5, 1pctCO2, esmFixClim1, and
18 esmFdbk1. It is the same set of models analyzed in Wenzel et al. (2016) and Winkler et al. (2019). The individual model setups
19 and components are illustrated in more detail in various studies, such as Arora et al. (2013); Wenzel et al. (2014); Mahowald
20 et al. (2016); Winkler et al. (2019).

21

22 The esmHistorical simulation spanned the period 1850 to 2005 and was driven by observed conditions such as solar forcing,
23 emissions or concentrations of short-lived species and natural and anthropogenic aerosols or their precursors, land use, anthro-
24 pogenic as well as volcanic influences on atmospheric composition. The models are forced by prescribed anthropogenic CO_2
25 emissions, rather than atmospheric CO_2 concentrations.

26

27 Several Representative Concentration Pathways (RCPs) have been formulated describing different trajectories of greenhouse
28 gas emissions, air pollutant production and land use changes for the 21st century. These scenarios have been designed based
29 on projections of human population growth, technological advancement and societal responses (van Vuuren et al., 2011; Tay-
30 lor et al., 2012). We analyzed simulations forced with specified concentrations of a high emissions scenario (RCP8.5) and
31 a medium mitigation scenario (RCP4.5) reaching a radiative forcing level of 8.5 and 4.5 W m^{-2} at the end of the century,
32 respectively. These simulations were initialized with the final state of the historical runs and spanned the period 2006 to 2100.

33

↓
at the end

1 1pctCO2 is an idealized fully coupled carbon-climate simulation initialized from a steady state of the pre-industrial control
2 run and atmospheric CO₂ concentration prescribed to increase 1% yr⁻¹ until quadrupling of the pre-industrial level. The sim-
3 ulations esmFixClim and esmFdbk aim to disentangle the two carbon cycle feedbacks in response to rising CO₂ analogous
4 to the 1pctCO2 setup: In esmFixClim CO₂-induced climate change is suppressed (i.e. radiation transfer model sees constant
5 pre-industrial CO₂ level), while the carbon cycle responds to increasing CO₂ concentration (*vice versa* for esmFdbk; Taylor
6 et al., 2009, 2012; Arora et al., 2013).

7

8 2.4 Estimation of greening sensitivities

9 We largely follow the methodology detailed in Winkler et al. (2019). For both model and observational data, the two-dimensional
10 global fields of LAI and the driver variables are cropped according to different classification schemes (namely, climatic regimes,
11 latitudinal bands and vegetation classes; Olson et al., 2001; Fritz et al., 2015). The aggregated values are area-weighted, aver-
12 aged in space, and temporally reduced to annual estimates dependent on the variable: annual maximum LAI, annual average
13 atmospheric CO₂ concentration, and growing degree days (GDD0, yearly accumulated temperature of days where near-surface
14 air temperature > 0° C).

15

16 We use a standard linear regression model to derive the historical greening sensitivities in models and observations alike (for
17 details see the Methods section *Estimation of historical LAI_{max} sensitivity* in Winkler et al., 2019). On the global scale, LAI_{max}
18 is assumed to be a linear function of atmospheric CO₂ concentration. For the temperature-limited high northern latitudes, we
19 also have to account for warming and include temperature as an additional driver. We do this using GDD0. Through a principal
20 component analysis (PCA) of CO₂ and GDD0 we avoid redundancy from co-linearity between the two driver variables, but
21 retain their underlying time-trend and interannual variability (for details see the Methods section *Dimension reduction using*
22 *principal component analysis* in Winkler et al., 2019). In particular, the PCA is performed on large-scale aggregated values
23 as well as on pixel level to investigate on spatial variations. We only retain the first principal component (denoted ω), which
24 explains a large fraction of the variance in models and observations (for more details see Supplementary Table 1 in Winkler
25 et al., 2019). Figure A1 depicts the temporal development of CO₂ and GDD0 as well as their principal component ω for
26 observations. For the NHL, LAI_{max} is then formulated as a linear function of the proxy driver time series ω (Winkler et al.,
27 2019). The best-fit gradients and associated standard errors of the linear regression model represent the LAI_{max} sensitivities,
28 or greening sensitivities, and their uncertainty estimates, respectively.

29

1 3 Results and Discussion

2 There are two parts to the EC methodology (Fig. 1) – a statistically robust relationship between modeled matching pairs of
3 predictor-predictand values and an observed value of the predictor. The predictors are from a representative historical period.
4 The predictands are modeled changes in a variable of interest at another forcing state of the system (e.g. potential future).
5 The projection of the observed predictor on the modeled relation yields a constrained value of the predictand. A causal basis
6 has to buttress the predictor-predictand relationship, else the EC method may be spurious. For example, meaningful coupling
7 between concurrent changes in GPP and LAI_{max} with increasing atmospheric CO₂ concentration underpins our specific case
8 study in the NHL, i.e. some of the enhanced GPP due to rising CO₂ concentration is invested in additional green leaves by
9 plants (Myneni et al., 1997a; Forkel et al., 2016; Zhu et al., 2016; Mao et al., 2016; Winkler et al., 2019). Supplementary Figure
10 1 in Winkler et al. (2019) illustrates the specifics of the causal link underlying this predictor-predictand relationship. This tight
11 coupling assures an approximately constant ratio of predictand to predictor across the models within the ensemble, thus setting
12 up the potential for deriving an EC estimate. Uncertainty ^{in/of} on the constrained estimate depends on the observed predictor and
13 modeled relationship, aside from the goodness-of-fit of the latter (Fig. 1). These are detailed below.

14

15 3.1 Uncertainty in Observed Predictor Due to Data Source

16 We investigate observational uncertainty using LAI data from two different sources, AVHRR (1/12 degree) and MODIS (1/20
17 degree), and spatially aggregating these ^{over} by broad vegetation classes, latitudinal bands and climatic regimes. The observed
18 large-scale LAI_{max} sensitivities to CO₂ forcing are always positive (greening), irrespective of the source data and the method
19 of aggregation (Fig. 2, Tab. 1). Overall, MODIS based estimates have higher uncertainty because of the shorter length of the
20 data record (17 years). The failure to reliably estimate sensitivities in tropical forests (also in the latitudinal band 30° S – 30°
21 N, and in hot, wet and humid climatic regimes, see Tab. 1 and Fig. 2) is due to saturation of optical remote sensing data over
22 dense vegetation (LAI_{max} > 5) and problems associated with high aerosol content and ubiquitous cloudiness. In other regions,
23 the estimated sensitivities are comparable across sensors and aggregation schemes, in particular in the high latitudinal band (>
24 60° N/S; AVHRR: $[3.4 \pm 0.5] \times 10^{-3}$, MODIS: $[3.6 \pm 0.9] \times 10^{-3} \text{ m}^2 \text{ m}^{-2} \text{ ppm}^{-1} \text{ CO}_2$). This aligns with previous studies
25 reporting a net increase in green leaf area across the high latitudes during the observational period (Myneni et al., 1997b; Zhu
26 et al., 2016; Forkel et al., 2016).

27

28 This analysis illustrates the applicability and limitations of using observed greening sensitivities to CO₂ forcing as a con-
29 straint on photosynthetic production. For example, data from both AVHRR and MODIS sensors provide a comparable estimate
30 of greening sensitivity in the colder high latitudes (boreal forests and tundra vegetation classes; Winkler et al., 2019). In the
31 lower latitudes, however, the discrepancies among the two sensors indicate a considerable observational uncertainty and thus
32 no robust estimation of the observed predictor is possible.

33

1 3.2 Uncertainty Due to Spatial Aggregation

2 We focus further analyses on the NHL region ($> 60^\circ \text{ N}$; Fig. 2b), because of two reasons. First, the direct human impact (i.e.
3 land management) can be neglected in the high latitudes, thus, we can assume that the observed changes reflect the response of
4 natural ecosystems. Second, the observational evidence of an increased plant productivity in the recent decades is well estab-
5 lished (e.g. Keeling et al., 1996; Myneni et al., 1997a; Graven et al., 2013; Forkel et al., 2016; Wenzel et al., 2016, and Sect.
6 3.1) – an important requisite in defining a robust predictor.

7

8 In addition to the physiological effect of CO_2 , also warming plays a key role in controlling plant productivity of the NHL
9 temperature-limited ecosystems, and thus, vegetation greenness. To avoid redundancy from co-linearity between CO_2 and
10 GDD0, we reduce dimensionality by performing a principal component analysis of the two driver variables (Sect. 2.4). The
11 resulting first principal component explains most of the variance and retains the trend and year-to-year fluctuations in both
12 CO_2 and GDD0. Therefore, we obtain a proxy driver (hereafter denoted ω) that represents the overall forcing signal causing
13 observed vegetation greenness changes in NHL (Fig. A1). Accordingly, greening sensitivity for the entire NHL area is derived
14 as response to ω , the combined forcing signal of rising CO_2 and warming. This procedure also enables a better comparability
15 between observations and models because varying strengths of physiological and radiative effects of CO_2 among models are
16 taken into account (Sect. 3.3 – 3.5).

17

18 The vegetated landscape in the NHL region is heterogeneous, with boreal forests in the south, vast tundra grasslands to the
19 north and shrublands in-between. The species within each of these broad vegetation classes respond differently to changes in
20 key environmental factors. Even within a species, such responses might vary due to different boundary conditions, such as
21 topography, soil fertility, micrometeorological conditions, etc. How this fine scale variation in greening sensitivity impacts the
22 aggregated value is assessed below.

23

24 The distribution of greening sensitivities from all NHL pixels is slightly skewed towards the positive (blue histogram). The
25 mean value of this distribution (blue dashed line) is comparable to the sensitivity estimate derived from the spatially-averaged
26 NHL time series (yellow dashed line; Fig. 3). Based on the Mann-Kendall test ($p > 0.1$), nearly over half the pixels (54%) show
27 positive statistically significant trends (greening), while about 10% show browning trends (possibly due to disturbances; Goetz
28 et al., 2005). The distribution of these statistically significant sensitivities (red histogram) therefore has two modes, a weak
29 browning and a dominant greening mode, resulting in a substantially higher mean value (red dashed line) in comparison to the
30 spatially-averaged estimate (yellow dashed line; Fig. 3). Thus, by taking into account the remaining 36% of non-significantly
31 changing pixels (as in the NHL spatially-averaged estimate), an additional source of uncertainty is possibly introduced. The
32 mean sensitivity value is, of course, higher when only pixels showing a greening trend are considered in the analysis (green
33 dashed line; Fig. 3). These are the only areas in NHL that actually show a large increase in plant productivity and consequently

1 significant changes in leaf area.

2

3 Model output of several ESMs (CMIP5) reveal similar pixel-level variation in both the predictor (LAI_{max} to ω , historical
4 simulation; Sect. 2.3) and associated changes in the predictand (GPP, 1pctCO₂; Sect. 2.3), although ESMs operate on much
5 coarser resolution (Fig. A2; see also Anav et al., 2013, 2015). Due to the coupling of the predictor and predictand, the distri-
6 bution of pixels with significant changes is approximately the same for the two variables (Fig. A2). Accordingly, averaging
7 the equally distributed estimates likely does not affect the predictor-predictand relationship in the model ensemble (Fig. 1).
8 Consequently, if all spatial gridded data arrays are consistently processed to spatially-aggregated estimates, each predictand
9 and predictor (observed and modeled) estimate contain a coherent component of spatial variations. In other words, considering
10 browning and non-significant pixels results in a lower overall LAI_{max} sensitivity in NHL, which in turn leads to a lower con-
11 strained estimate of ΔGPP in NHL. This is consistent with the underlying relationship between predictor and predictand. On a
12 related note, Bracegirdle and Stephenson (2012a) suggest that this source of error is not significantly dependent on the spatial
13 resolution when comparing model subsets from high to low resolution.

14

15 The above analysis informs that spatially-averaged estimates are approximations containing a random error component due
16 to inclusion of data from insignificantly changing pixels and a systematic bias component from pixels of reversed sign. This
17 uncertainty is relevant to the EC method, where the observed sensitivity decisively determines the constrained estimate from
18 the ensemble of ESM projections (Kwiatkowski et al., 2017; Winkler et al., 2019). However, if spatial variations are treated
19 consistently as an inherent component of observations and models, the EC method is only slightly susceptible to this source of
20 uncertainty.

21

22 3.3 Uncertainty Due to Temporal Variations

23 We seek recourse to longterm CMIP5 ESM simulations covering the historical period 1850 to 2005 (Sect. 2.3) to assess
24 temporal variation in the predictor variable, because of the shortness of observational record. Three representative models
25 (CESM1-BGC, MIROC-ESM, and HadGEM2-ES) spanning the full range of NHL greening sensitivities in the CMIP5 en-
26 semble (Winkler et al., 2019) are selected for this analysis. For each model, LAI_{max} sensitivity to ω in moving windows of
27 different lengths are evaluated (15, 30, and 45 years; Fig. 4 and A3). The analysis reveals two crucial aspects that highlight how
28 temporal variations impair comparability of the predictor variable between models and observations – an essential component
29 of the EC approach.

30

31 First, window locations of modeled and observed predictor variable have to match. If the forcing in the simulations is low,
32 for example, as in the second half of the 19th century when CO₂ concentration was increasing slowly, inter-annual variability
33 dominates and LAI_{max} sensitivity cannot be accurately estimated irrespective of the window length (Fig. 4 and A3). With
34 increasing forcing over time (rising yearly rate of CO₂ infusion, and consequently, the concentration), the signal-to-noise ratio

?

1 increases and LAI_{max} sensitivity to ω estimation stabilizes, for example, as in the second half of the 20th century. Therefore,
2 LAI_{max} sensitivities estimated at different temporal locations result in non-comparable values and eventually a false con-
3 strained estimate (details in Sect. 3.4). As an example, modeled sensitivities based on a 30-year window centered on year 1900,
4 when CO_2 level increased by 10 ppm, and observed sensitivity estimated from a 30-year window centered on year 2000, when
5 CO_2 level increased by 55 ppm, describe different states of the system and therefore should not be contrasted in the EC method.

6
7 Second, in addition to temporal location, also window lengths have to match between observations and models. For all three
8 models, sensitivities estimated from 15-year chunks show high variability and thus, a 15-year record is perhaps too short to
9 obtain robust estimates. The LAI_{max} sensitivity estimation becomes more stable with strengthening forcing and increasing
10 window length (Fig. 4 and A3). As a consequence, using short-term observed sensitivity as a constraint on long-term model
11 projections results in an incorrect EC estimate. Hence, the MODIS sensor record is, on the one hand, too short and does not,
12 on the other hand, overlap temporally with the historical CMIP5 forcing. Therefore, it does not provide a robust predictor in
13 this EC study.

15 3.4 Level and Time Rate of CO_2 Forcing

16 The EC method raises an obvious question – does it not implicitly assume that the key operative mechanisms underpinning the
17 EC relation remain unchanged because a future system state is being predicted based on its past behavior? To be specific, we
18 are attempting to predict GPP at a future point in time based on greening sensitivity inferred from the past. Does this not require
19 the assumption that the key underlying relationship which makes this prediction possible, namely, a robust coupling between
20 contemporaneous changes in GPP and LAI_{max} remains unchanged from the past to the future? To address this question, we
21 resort to the CMIP5 idealized simulation (1pctCO2), where atmospheric CO_2 concentration increases 1% annually, starting
22 from a pre-industrial level of 284 ppm until a quadruple of this value is reached (Sect. 2.3). We limit the analysis to the three
23 models (CESM1-BGC, MIROC-ESM, and HadGEM2-ES) which bracket the full range of GPP enhancement and LAI_{max}
24 sensitivity in the original seven ESM ensemble (Winkler et al., 2019).

25
26 The relationship between simultaneous changes in GPP and LAI_{max} remains linear for all CMIP5 models in the range
27 $1 \times CO_2$ to $2 \times CO_2$ (Fig. 5 and A4, Tab. 2). With concentration increasing beyond $2 \times CO_2$, all models show weakening correla-
28 tion (R^2 , Tab. 2) and decreasing slope (b , Tab. 2) of this relationship (Fig. 5 and A4), suggesting a saturating rate of allocation
29 of additional GPP to new leaves at higher levels of CO_2 . Consequently, LAI_{max} sensitivity to increasing CO_2 and associated
30 warming decreases. At and over $4 \times CO_2$ (1140 ppm), a level unlikely to be seen in the near future, there appears to be no
31 relationship between ΔGPP and ΔLAI_{max} . This raises the question as to what extent does the weakening of the relationship
32 between the predictor and predictand in each model at higher CO_2 concentrations affect the EC analysis (Fig. 1). To shed light
33 on this matter, we perform the following thought experiment.

34
in some
models

1 Understanding the relationship and interplay between forcing (increasing CO₂ concentration), predictor (LAI_{max} sensitiv-
 2 ity), and the predictand (Δ GPP) is key to evaluating the EC method. We conceive four possible scenarios of how the sys-
 3 tem might behave with increasing forcing. For simplicity, we assume linearly increasing CO₂ concentration, LAI represents
 4 LAI_{max}, and GPP refers to its annual value below (Fig. 6). The four scenarios are: *All linear*, *all non-linear* (saturation), and
 5 two *mixed linear / non-linear* cases (Tab. A1). We emulate a multi-model ensemble by applying different random parameteri-
 6 zations for the linear and saturation (the hyperbolic tangent function) responses of GPP to CO₂ and of LAI to GPP. One of these
 7 realizations is assumed to represent pseudo-observations (dashed lines, Fig. 6). We discuss one case in detail for illustrative
 8 purposes (No. 3, Tab. A1).

9

10 In scenario 3, Δ GPP increases linearly with increasing CO₂ (Fig. 6a), while Δ LAI/ Δ GPP saturates (Fig. 6b). The LAI sen-
 11 sitivity to CO₂ weakens with increasing forcing (Fig. 6c) as a response to saturation of GPP allocation to leaf area. We derive
 12 LAI sensitivities to CO₂ for three different periods ('past periods' in Fig. 6c) to constrain Δ GPP at a much higher CO₂ level
 13 ('projected period' in Fig. 6a). Next, we apply the EC method on these pseudo-projections of Δ GPP relying on LAI sensitivi-
 14 ties derived from the three past periods (Fig. 6d). The EC method is applicable even at a low forcing level (past period 1) in this
 15 simplified scenario because we neglect stochastic internal variability of the system. The slope of emergent linear relationship
 16 increases (Fig. 6d) as modeled LAI sensitivities decrease with rising CO₂ concentration (Fig. 6c). The observational constraint
 17 on future Δ GPP, however, remains nearly the same, because pseudo-observed LAI sensitivity also weakens at higher CO₂
 18 levels (dashed lines, Fig. 6c, d). Thus, the three EC estimates of Δ GPP are approximately identical (Fig. 6d) and independent
 19 of the forcing level during past periods. With intensified forcing, the relationship between predictor and predictand remains
 20 linear within the model ensemble, although their relationship becomes non-linear within each model and, crucially, in reality
 21 as well. In other words, as long as the models agree on the occurrence and strength of saturation for given forcing, i.e. the
 22 dynamics of the system, the inter-model variations of predictor and predictand relate linearly within the ensemble (Fig. 6). The
 23 same behavior is also seen in the other three scenarios (Tab. A1; Fig. A5, A6).

24

25 Nevertheless, with ever increasing forcing and associated steepening of the emergent linear relationship, the LAI sensitivity
 26 loses its explanatory power at some point because the linear relationship eventually lies within the observational uncertainty
 27 and no meaningful constraint can be derived. This and disagreement between models on system dynamics are ultimate limits
 28 of the EC method. Interestingly, we find that all CMIP5 models agree on the occurrence of saturation, but slightly disagree on
 29 the strength of saturation for given CO₂ forcing (Fig. 5, A4, and Tab. 2). Further, we find that the 'all non-linear' scenario best
 30 describes the dynamics of the system in the forcing range from 1×CO₂ to 4×CO₂. However, the saturation of LAI to GPP
 31 happens at a lower CO₂ level than saturation of GPP to CO₂. Still, inferences from interpretation of Case 3 (Fig. 6) are equally
 32 applicable.

33

34 Results from the above thought experiment also highlight the importance of matching window locations and lengths between
 35 models and observations, as discussed earlier (Sect. 3.3). For instance, taking LAI sensitivity from past period 2 (green dashed

1 line, Fig. 6d) as an observational constraint on the multi-model linear relationship based on past period 3 (red solid line, Fig.
2 6d), results in a significant overestimation of constrained ΔGPP (intersection of the two lines, Fig. 6d).

3

4 The above analysis informs that the constrained GPP estimate at one future period (e.g. $2\times\text{CO}_2$) is nearly independent of
5 the past periods from when the observational sensitivities are derived, for most realistic scenarios. Now, we evaluate the EC
6 method where sensitivity from one past period is used to obtain constrained GPP estimates at different periods in a potential
7 future, i.e. progressively farther down the time-line of a CO_2 -enriched world. We utilize the greening sensitivity derived from
8 35 years of observed LAI_{max} data (AVHRR, Sect. 2.1) and apply the EC method to CMIP5 1pctCO2 simulations. The sensi-
9 tivities in this case are due to forcing from both CO_2 increase and associated warming during the observational period (Sect.
10 2.4). We seek constrained GPP estimates for the NHL at different CO_2 levels ($2\times\text{CO}_2$, $3\times\text{CO}_2$, and $4\times\text{CO}_2$).

11

12 Winkler et al. (2019) previously reported a strong linear relationship between modeled contemporaneous changes in LAI_{max}
13 and GPP arising from the combined radiative and physiological effects of CO_2 enrichment until $2\times\text{CO}_2$ in the CMIP5 ensem-
14 ble. As a result, models with low LAI_{max} sensitivity to ω project lower ΔGPP for a given increment of CO_2 concentration, and
15 *vice versa*. Thus, the large variation in modeled historical LAI_{max} sensitivities linearly maps to variation in ΔGPP at $2\times\text{CO}_2$
16 (Winkler et al., 2019, blue line, Fig. 7a). At higher levels, such as $3\times\text{CO}_2$ (green line, $R^2 = 0.93$) and $4\times\text{CO}_2$ (red line, R^2
17 $= 0.88$), this linear relationship within the model ensemble, while still present, weakens (Fig. 7a; Tab. 3). This is because the
18 CMIP5 models do not agree on the strength of the saturation effect at higher CO_2 levels (Fig. 5 and A4). The increment in
19 constrained GPP estimates for successive equal increments of CO_2 decreases due to the saturation effect in all CMIP5 models
20 (dashed horizontal lines, Fig. 7a). For example, the change in GPP between $3\times\text{CO}_2$ and $4\times\text{CO}_2$ ($\Delta\text{GPP} \sim 1.06 \text{ Pg C yr}^{-1}$,
21 Tab. 3) is much lower than between $2\times\text{CO}_2$ and $3\times\text{CO}_2$ ($\Delta\text{GPP} \sim 2.34 \text{ Pg C yr}^{-1}$, Tab. 3).

22

23 We have thus far focused on the magnitude of CO_2 concentration change and not on the time rate of this change. For example,
24 a given amount of change in CO_2 concentration, say 200 ppm, can be realized over different time periods, say over a 100 or 150
25 years. The problem of varying rates of CO_2 concentration change is implicitly encountered when ESMs are executed under
26 different forcing scenarios, such as RCPs (Sect. 2.3). A question then arises whether the constrained predictand estimate is
27 independent of the time rate of CO_2 concentration change and dependent only on the magnitude of CO_2 concentration change.
28 To investigate this aspect of forcing, we extract GPP estimates at the same CO_2 concentration (535 ppm; final concentration
29 in RCP4.5) from three simulations of different forcing rates and calculate the difference relative to a common initial CO_2
30 concentration (380 ppm; initial concentration of RCP scenarios). Hence, the magnitude of the forcing is the same but applied
31 over different durations (RCP4.5: $\sim 90\text{yr}$, RCP8.5: $\sim 45\text{yr}$, and 1pctCO2: $\sim 30\text{yr}$). A clear majority of the CMIP5 models show
32 substantial differences in ΔGPP between the different pathways of CO_2 forcing. In general, GPP changes are higher for lower
33 time rates of CO_2 forcing, i.e. forcing over longer time periods. As a consequence, the EC estimates of ΔGPP for the same
34 increase in CO_2 concentration are scenario-dependent (Fig. 7b; Tab. 3) – a counter-intuitive result. For instance, ΔGPP in the
35 low- CO_2 -rate scenario (RCP4.5: $\Delta\text{GPP} \sim 2.84 \text{ Pg C yr}^{-1}$, Tab. 3) is $\sim 39\%$ (1pctCO2: $\Delta\text{GPP} \sim 2.05 \text{ Pg C yr}^{-1}$, Tab. 3) and

is this the 1pctCO2 scenario

1 ~20% (RCP8.5: $\Delta GPP \sim 2.38 \text{ Pg C yr}^{-1}$, Tab. 3) higher than the high-CO₂-rate scenarios for an increase of 155 ppm CO₂.
2 This analysis suggests that the vegetation response to rising CO₂ is pathway dependent, at least in the NHL. One of the reasons
3 for this could be species compositional changes in scenarios of low forcing rates, i.e. over longer time frames. This novel result,
4 however, requires a separate in-depth study.

5 3.5 Effects of CO₂ Forcing

6 Higher concentration of CO₂ in the atmosphere stimulates plant productivity through the fertilization and radiative effects (Ne-
7 mani et al., 2003; Leakey et al., 2009; Arora et al., 2011; Goll et al., 2017). The two effects can be disentangled in the model
8 world by conducting simulations in a 'CO₂ fertilization effect only' (esmFixClim1) and a 'radiative effect only' (esmFdbk1)
9 setup (Sect. 2.3). These are termed below as idealized model simulations. We investigate here whether historical runs and
10 observations, which include both effects, can be used to constrain GPP changes in idealized CMIP5 simulations (e.g. as in
11 Wenzel et al., 2016).

12
13 We find strong linear relationships between historical LAI_{max} sensitivity and ΔGPP for $2 \times CO_2$ in both idealized setups
14 (esmFixClim1: $R^2 = 0.92$, esmFdbk1: $R^2 = 0.98$, Tab. 3, Fig. 7c). Consequently, this linear relationship is also pronounced for
15 calculated sums of both effects for each model (esmFixClim1 + esmFdbk1: $R^2 = 0.95$, Tab. 3, Fig. 7c). This suggests that the
16 two effects act additively on plant productivity and, thus, each effect can be simply expressed in terms of a scaling factor of
17 the total GPP enhancement. Hence, the application of the EC method on idealized simulations using real world observations is
18 conceptually feasible.

19
20 Interestingly, the two effects contribute about the same to the general increase in GPP at $2 \times CO_2$ (esmFixClim1: ΔGPP
21 $\sim 1.35 \text{ Pg C yr}^{-1}$, esmFdbk1: $\Delta GPP \sim 1.38 \text{ Pg C yr}^{-1}$, Tab. 3, Fig. 7c). At higher concentrations, such as $3 \times CO_2$ and $4 \times CO_2$,
22 the enhancement in GPP saturates in both idealized setups. However, the radiative effect becomes dominant relative to the
23 CO₂ fertilization effect when CO₂ concentration exceeds $2 \times CO_2$ (e.g. at $4 \times CO_2$ esmFixClim1: $\Delta GPP \sim 2.42 \text{ Pg C yr}^{-1}$,
24 esmFdbk1: $\Delta GPP \sim 3.06 \text{ Pg C yr}^{-1}$, Tab. 3). Therefore, we can expect that at some point in the future, NHL photosynthetic
25 carbon fixation will benefit more from climate change (e.g. warming) than from the fertilizing effect of CO₂.

A

26 3.6 Uncertainties in the Multi-Model Ensemble

27 Besides methodological sources of uncertainty discussed above, the estimate of an EC may also be deficient due to inaccurate
28 assumptions about the model ensemble. First, possible common systematic errors in a multi-model ensemble (i.e. the entire
29 ensemble misses an unknown but for the future essential process) are implicitly omitted in the EC approach, however, could
30 cause a general over- or underestimation of the constrained value (Bracegirdle and Stephenson, 2012b; Stephenson et al.,
31 2012). Second, the set of forcing variables for historical simulations may be incomplete (i.e. not yet identified drivers of
32 observed changes) and thus the comparability of observations and model simulations is limited (Flato et al., 2013). Third,
33 the EC method can be overly sensitive to individual models of the ensemble, which has a bearing on the robustness of the

reword

1 constrained value (Bracegirdle and Stephenson, 2012b). Bracegirdle and Stephenson (2012b) proposed a diagnostic metric
2 (Cook's distance) to test an ensemble for influential models. Fourth, the predictand-predictor relationship not only has to rely
3 on a physical, but also on a logical connection within the model ensemble. For instance, Wenzel et al. (2016) established a
4 linear relationship between relative changes in the predictand taking the initial state into account (changes in GPP for doubling
5 of CO₂ relative to the initial pre-industrial state), and a predictor neglecting the initial state (historical sensitivity of CO₂
6 amplitude to rising CO₂). This statistical relationship can be spurious, because the model skill of simulating an accurate initial
7 state and a plausible sensitivity to a forcing are not connected. These issues are to be contemplated when establishing an EC
8 estimate and evaluating its robustness.

9 **4 Conclusions**

10 An in-depth analysis of the EC method is illustrated in this article through its application to projections of change in NHL
11 photosynthesis under conditions of rising atmospheric CO₂ concentration. Key conclusions highlighting the functionality of
12 the EC method are presented below.

13
14 The importance of how the observational predictor is obtained cannot be emphasized enough because the EC method is
15 particularly sensitive to observational uncertainty. The single observational estimate essentially determines the EC, whereas
16 the emergent linear relationship is established based on a collection of multi-model estimates (each model gets 'one vote',
17 however, some models might be more influential than others; Bracegirdle and Stephenson, 2012b). Hence, the observational
18 uncertainty has a much larger bearing on the EC than the uncertainty of each individual model. To overcome this source of
19 uncertainty, various meaningful observations should be taken into consideration when establishing the observed predictor.

20
21 Spatially aggregating observations and model output of different resolutions in the EC method constitutes another source
22 of uncertainty. Predictors and predictands expressed as regional estimates (e.g. area-weighted mean of the NHL) are approxi-
23 mations of complex fine-scale processes. Aggregation will inevitably introduce a random error component due to inclusion of
24 estimates from areas where the predictor is not changing or a systematic bias from areas where the predictor has a reversed
25 sign. Thus, the spatially-aggregated variables are meaningful only if most of the region is in agreement about the response to
26 CO₂ forcing (e.g. more than half of the NHL is greening with rising CO₂). However, we find that the source of uncertainty
27 related to spatial aggregation is of minor importance as long as spatial variations in observations and models simulations are
28 treated consistently.

29
30 A large source of uncertainty is associated with temporal variability of the predictor variable when comparing models and
31 observations. Establishing a robust predictor requires evaluating temporal window lengths of sufficient duration (approximately
32 30 years) and their locations along the forcing time line. Both window length and location should match between models and
33 observations in the EC method. For example, the analysis in Wenzel et al. (2016) might have yielded different results and

1 conclusions if model and observational predictor sensitivities were temporally matched. We find that the relevance of window
2 length decreases with increasing and accelerating forcing, depending on the magnitude of natural/internal variability (signal-
3 to-noise ratio) of the predictor variable.

4
5 The level, effect and time-rate of applied CO₂ forcing can have a bearing on the linear relationship between the predictand
6 and predictor variables (Fig. 1). In our case study, the relationship underpinning the EC method, namely, that between concur-
7 rent ΔGPP and ΔLAI_{max} changes non-linearly with increasing forcing level (i.e. saturation with rising CO₂ concentration).
8 The EC method can still be applied, because the CMIP5 models agree on the non-linear behavior of the system. However,
9 at very high CO₂ concentrations the models diverge and this relation breaks down, at which point the EC method fails. The
10 two dominant effects of rising CO₂ concentration on vegetation, namely, the fertilization and radiative effects, appear to be
11 approximately additive in terms of GPP enhancement to CO₂ forcing in the NHL. Therefore, the EC method can be applied
12 to constrain estimates of GPP due to one or the other, or both the effects. The models, however, document a higher radiative
13 effect than fertilization at high CO₂ concentrations, i.e. 3×CO₂ and higher. Another intriguing conclusion from our analysis
14 is that the time-rate of forcing has an effect on GPP changes, that is, the projected GPP enhancement to CO₂ forcing seems
15 to be dependent on how the forcing is applied over time, as in different scenarios or RCPs. This aspect is presently not well
16 understood and requires further study.

17
18 The EC framework is widely promoted as observation-based evaluation tool for climate projections, especially in the context
19 of the nascent CMIP6 ensemble (Eyring et al., 2019; Hall et al., 2019). Previous EC studies, however, exclusively focused on
20 predictor-predictand combinations which exhibit so-called existent ECs (Hall et al., 2019), i.e. predictor and predictand are
21 found to relate linearly across the ensemble. In the context of ESM evaluation, non-existent ECs, i.e. predictor and predictand
22 are found to be unrelated in the ensemble, are equally important. Since predictor and predictand variables are premised on
23 our mechanistic process understanding, non-existent ECs reveal a fundamental disagreement on the system dynamics among
24 the models. This study encourages to scrutinize these system dynamics in the predictor-predictand space and also report such
25 non-existent, yet expected, ECs in order to advance model development and evaluation.

26
27 Across different disciplines each EC and its set of predictor and predictand are unique to some extent and require an individ-
28 ual detailed examination. In this article, we addressed general potential sources of uncertainty and limitations in the EC method
29 by the means of a case study in carbon cycle research. ~~Thus,~~ the illustrated results are qualitatively transmissive to other sets
30 of predictors and predictands and are generally relevant in Earth system sciences.

Make (A) & (B) consistent.

1 *Author contributions.* A.J.W. performed the analysis. All authors contributed ideas and to writing of the manuscript.

2 *Competing interests.* The authors declare that they have no conflict of interest.

3 *Acknowledgements.* We thankfully acknowledge T. Park and C. Chen for their help with remote sensing data. We thank G. Lasslop for
4 reviewing the manuscript. R.B.M. thanks Alexander von Humboldt Foundation and NASA's Earth Science Division for funding support that
5 made his participation possible in this research.

Please thank the reviewers

Were these values determined using grid values for NHL.

1 **Table 2.** Slopes (b) and coefficients of determination (R^2) for regression between changes of LAI_{max} against changes in annual mean GPP
 2 at different atmospheric CO_2 levels in all available CMIP5 models (1pctCO2 simulation). Asterisks denote non-significant values: ** $p >$
 3 0.1; * $p > 0.05$.

Correlation details	< 2xCO ₂		> 2xCO ₂ & < 3xCO ₂		> 3xCO ₂	
	b	R^2	b	R^2	b	R^2
MIROC-ESM	0.23	0.97	0.16	0.89	0.08	0.63
CESM1-BGC	0.45	0.93	0.36	0.82	0.27	0.62
4 GFDL-ESM2M	0.37	0.89	0.04	0.07**	0.01	0.12**
CanESM2	0.22	0.95	0.19	0.83	0.17	0.67
HadGEM2-ES	0.13	0.99	0.08	0.96	0.06	0.78
MPI-ESM-LR	0.13	0.94	0.09	0.78	0.04	0.51
NorESM1-ME	0.26	0.94	0.2	0.77	0.09	0.27